

# Revista de Estudios Marítimos y Sociales

*Publicación científica de carácter semestral*

Año 17 - Número 24 - ene-jun de 2024 - Mar del Plata - Argentina - ISSN 2545-6237

## Cómo se incorporan (o no) los eventos en los sets de datos: Los obstáculos y las orientaciones en la utilización de los periódicos en el estudio de los conflictos

*Translation: How Events Enter (or Not) Data Sets: The Pitfalls and Guidelines of Using Newspapers in the Study of Conflict*

ARK CAICYT: <http://id.caicyt.gov.ar/ark:/s25456237/z4uar6b01>

Leila Demarest\*

Instituto de Ciencias Políticas, Universidad de Leiden

Arnim Langer<sup>^</sup>

Centro de Investigación sobre la Paz y el Desarrollo. Universidad de Lovaina

Traducción: Daniel Nicolás Rabino •

Instituto de Humanidades y Ciencias Sociales – Consejo Nacional de Investigaciones Científicas y Técnicas

Revisión del inglés: Marlene González Marín\*

Universidad Nacional de Mar del Plata, Facultad de Humanidades. Departamento de Lenguas Modernas.

\* Profesora adjunta de Política Africana. Leiden, Países Bajos. Correo electrónico: [l.demarest@fsw.leidenuniv.nl](mailto:l.demarest@fsw.leidenuniv.nl)  
ORCID 0000-0001-6887-9937

<sup>^</sup> Profesor de Relaciones Internacionales en la Universidad Católica de Lovaina, Bélgica.

• Doctorando en Historia. Profesor y Licenciado en Historia. Mar del Plata, Buenos Aires, Argentina. Correo electrónico: [nicolas236@gmail.com](mailto:nicolas236@gmail.com) ORCID 0000-0002-1789-6435.

<sup>v</sup> Profesora de Inglés. Mar del Plata, Buenos Aires, Argentina Correo electrónico: [marlugm@yahoo.com.ar](mailto:marlugm@yahoo.com.ar)



## Cómo se incorporan (o no) los eventos en los sets de datos: Los obstáculos y las orientaciones en la utilización de los periódicos en el estudio de los conflictos <sup>1</sup>

Translation

*How Events Enter (or Not) Data Sets: The Pitfalls and Guidelines of Using Newspapers in the Study of Conflict* <sup>2 3 4</sup>

Leila Demarest\* y Arnim Langer^

Traducción: Daniel Nicolás Rabino ♦

Revisión del Inglés: Marlene Gonzáles Marín ▽

Recibido: 14 de diciembre de 2023

Aceptado: 4 de enero de 2024

### Abstract

While conflict event data sets are increasingly used in contemporary conflict research, important concerns persist regarding the quality of the collected data. Such concerns are not necessarily new. Yet, because the methodological debate and evidence on potential errors remains scattered across different subdisciplines of social sciences, there is little consensus concerning proper reporting practices in codebooks, how best to deal with the different types of errors, and which

---

<sup>1</sup> Tomado de Leila Demarest y Arnim Langer, «How Events Enter (or Not) Data Sets: The Pitfalls and Guidelines of Using Newspapers in the Study of Conflict», *Sociological Methods & Research* 51, n.º 2 (1 de mayo de 2022): 632-66, <https://doi.org/10.1177/0049124119882453>. Agradecemos la autorización otorgada por los autores y a SAGE para traducir y publicar el presente texto en nuestra revista.

<sup>2</sup> Enlace al artículo original <https://journals.sagepub.com/doi/epub/10.1177/0049124119882453>

<sup>3</sup> <https://creativecommons.org/licenses/by-nc/4.0/legalcode>

<sup>4</sup> <https://creativecommons.org/licenses/by-nc/4.0/>

\* Instituto de Ciencias Políticas, Universidad de Leiden, edificio Pieter de la Court, Wassenaarseweg 52, 2333 AK Leiden, Países Bajos. Correo electrónico: [l.demarest@fsw.leidenuniv.nl](mailto:l.demarest@fsw.leidenuniv.nl) ORCID 0000-0001-6887-9937

^ Centro de Investigación sobre la Paz y el Desarrollo (CRPD), Universidad de Lovaina, Bélgica

♦ Doctorando en Historia. Profesor y Licenciado en Historia por la Universidad Nacional de Mar del Plata. Mar del Plata, Buenos Aires, Argentina. ORCID 0000-0002-1789-6435. Correo electrónico: [nicolas236@gmail.com](mailto:nicolas236@gmail.com)

▽ Profesora de Inglés. Mar del Plata, Buenos Aires, Argentina Correo electrónico: [marlugm@yahoo.com.ar](mailto:marlugm@yahoo.com.ar)



types of errors should be prioritised. In this article, we introduce a new analytical framework—that is, the Total Event Error (TEE) framework—which aims to elucidate the methodological challenges and errors that may affect whether and how events are entered into conflict event data sets, drawing on different fields of study. Potential errors are diverse and may range from errors arising from the rationale of the media source (e.g., selection of certain types of events into the news) to errors occurring during the data collection process or the analysis phase. Based on the TEE framework, we propose a set of strategies to mitigate errors associated with the construction and use of conflict event data sets. We also identify a number of important avenues for future research concerning the methodology of creating conflict event data sets.

**Keywords:** conflict events - data error - bias, -unreliability - media data - total survey error

### Resumen

Si bien los conjuntos de datos sobre acontecimientos conflictivos se utilizan cada vez más en la investigación de conflictos, persisten importantes preocupaciones sobre la calidad de los datos recogidos. Estas preocupaciones no son necesariamente nuevas. Sin embargo, debido a que el debate metodológico y las pruebas sobre los posibles errores siguen estando dispersos entre las distintas subdisciplinas de las ciencias sociales, hay poco consenso en cuanto a las prácticas adecuadas de notificación en los libros de códigos, la mejor manera de tratar los diferentes tipos de errores, y a qué tipos de errores hay que dar prioridad. En este artículo, introducimos un nuevo marco analítico, el entorno llamado en inglés Total Event Error (TEE), que tiene como objetivo dilucidar los retos metodológicos y los errores que pueden afectar a la introducción de los sucesos en los conjuntos de datos de conflictos y a la forma en que se hace, basándose en diferentes campos de estudio. Los errores potenciales son diversos y pueden ir desde errores derivados de la lógica de los medios de comunicación (por ejemplo, la selección de determinados tipos de sucesos en las noticias) a errores que se producen durante el proceso de recogida de datos o la fase de análisis. Basándonos en el entorno TEE, proponemos un conjunto de estrategias para mitigar los errores asociados con la construcción y el uso de conjuntos de datos de eventos conflictivos. También identificamos una serie de vías importantes para la investigación futura en relación con la metodología de creación de conjuntos de datos de eventos conflictivos.

**Palabras clave:** eventos conflictivos - error de datos – sesgo - falta de fiabilidad - datos de los medios de comunicación - error total de la encuesta



## Introducción

Con el giro cuantitativo en los estudios sobre la paz y los conflictos (por ejemplo, Collier y Hoeffler 2002, Fearon y Laitin 2003, Hirshleifer 1994:3), la recopilación de datos detallados sobre los eventos de conflicto, los actores y el número de víctimas han impulsado muchos proyectos de investigación en este campo. El desarrollo de conjuntos de datos sobre conflictos también se ha acelerado en los últimos años. Las principales tendencias en los nuevos proyectos de datos incluyen un mayor enfoque en los eventos de conflicto desagregados -tanto en el tiempo como en el espacio-, así como un enfoque en las formas de conflicto de bajo nivel, como las protestas, en contraposición a la guerra civil (Bernauer y Gleditsch 2012:375-78, Gleditsch et al. 2014:303-5, 308-9). Los informes de noticias han sido la fuente más importante de datos sobre los eventos conflictivos, ya que están ampliamente disponibles y a menudo son accesibles a bajo coste. Sin embargo, el uso generalizado de los informes de los medios de comunicación como fuente empírica plantea importantes preocupaciones sobre la calidad de los datos de los acontecimientos.

La recopilación de datos sobre acontecimientos políticos tiene una larga historia, tanto en la literatura sobre movimientos sociales (por ejemplo, Eisinger 1973) como en la de relaciones internacionales (por ejemplo, Azar 1980, McClelland 1976). La preocupación por la validez y la fiabilidad de los datos de los medios de comunicación para la investigación en ciencias políticas no es, por tanto, necesariamente nueva (por ejemplo, Danzger 1975, Franzosi 1987). No obstante, la mayor disponibilidad de fuentes mediáticas (en línea) y el desarrollo de nuevos conjuntos de datos han estimulado nuevos debates en este campo. Una característica importante de muchos de los nuevos conjuntos de datos es su enfoque geográfico en el mundo en desarrollo. Los ejemplos incluyen el conjunto de datos de localización y eventos de conflictos armados (ACLED, Raleigh et al. 2010), la base de datos de conflictos sociales en África (SCAD, Salehyan et al. 2012), el conjunto de datos de disturbios sociales urbanos en África y Asia (USDAA) (Urdal 2008), el conjunto de datos de eventos georreferenciados UCDP (UCDP GED, Sundberg y Melander 2013), la base de datos de terrorismo mundial (LaFree y Dugan 2007), el conjunto de datos de atrocidades mundiales del Grupo de Trabajo sobre Inestabilidad Política (PITF) (Schrodt y Ulfelder 2016), la base de datos

Demarest, Leila y Langer, Arnim "Cómo se incorporan (o no) los eventos en los sets de datos: Los obstáculos y las orientaciones en la utilización de los periódicos en el estudio de los conflictos" *Revista de Estudios Marítimos y Sociales*, N°24, ene-jun 2024, pp. 197-244.



de movilización de masas en las autocracias (Weidmann y Rød 2015), y el conjunto de datos de violencia unilateral de Konstanz (Schneider y Bussmann 2013). Por el contrario, gran parte del debate metodológico y las pruebas con respecto al uso de los informes de los medios de comunicación para construir datos de eventos se encuentran en la investigación de los movimientos sociales centrada en Occidente (por ejemplo, Earl et al. 2004, Hutter 2014), así como en los estudios de comunicación (Galtung y Ruge 1965, Harcup y O'Neill 2016, Krippendorff 2013). Además, aunque en los últimos años se han evaluado críticamente diferentes retos metodológicos asociados a la generación de nuevos conjuntos de datos de eventos conflictivos (por ejemplo, Eck 2012, Salehyan 2015, Weidmann 2015, 2016), este campo de estudio se beneficiaría de una mayor sistematización de los resultados de la investigación.

Para ello, el presente artículo introduce un nuevo marco analítico que capta los retos metodológicos de la utilización de los informes de noticias para generar datos sobre conflictos y reconoce una amplia gama de errores que pueden afectar a la introducción de eventos en los conjuntos de datos y a la forma de hacerlo. El entorno (framework) del Error Total de Eventos (TEE siglas en inglés) se basa en la literatura de investigación de encuestas y en el marco del Error Total de la Encuesta (TSE siglas en inglés). Por analogía con Groves y sus colegas (2004:41-63), distinguimos entre errores de medición y errores de representación. El entorno abarca formas bien conocidas de error mencionadas en la literatura, como el sesgo de selección, que surge cuando los periódicos seleccionan deliberadamente algunos eventos para su publicación, mientras que dejan otros eventos sin reportar (Earl et al. 2004:68-72, Jenkins y Maher 2016, Ortiz et al. 2005, Weidmann 2016). Sin embargo, también consideramos errores que no son necesariamente causados por la lógica de las fuentes mediáticas y que han recibido mucha menos atención en la literatura. Estos errores surgen durante el proceso de recogida de datos, como la codificación de las variables clave, o en la fase de análisis, cuando los investigadores recurren a valores imputados para los datos que faltan (por ejemplo, el lugar de un evento). Además, si bien el sesgo, o una diferencia sistemática entre el valor medido y el valor real, es una forma importante de error, también dirigimos la atención hacia la falta de fiabilidad o la desviación aleatoria del valor real, que socava la precisión.



El entorno de TEE ofrece un puente entre las ideas metodológicas de los estudios de los acontecimientos conflictivos y los estudios de los movimientos sociales y la comunicación centrados en Occidente. Esto tiene la ventaja de que ideas, métodos y procedimientos que son comunes en estas últimas literaturas se introducen y discuten en relación con los conflictos en los países en desarrollo. Prestamos especial atención a las implicaciones de centrarnos en los conflictos de los países en desarrollo en contraposición a los contextos occidentales. En efecto, mientras que los estudios centrados en Occidente se han centrado usualmente en los acontecimientos de protesta, nosotros nos centramos en una gama más amplia de acontecimientos, incluyendo las protestas, pero también los conflictos armados violentos. Por último, discutimos y comparamos las formas humanas y automatizadas de recopilación y codificación de datos. Aunque a menudo se ha expresado optimismo con respecto a las posibles oportunidades y ventajas de la codificación automatizada (por ejemplo, Bond et al. 1997, King y Lowe 2003, Schrodtt y van Brackle 2013), hasta el momento, no es ampliamente utilizado en los estudios de conflictos. De manera ilustrativa, la codificación humana es utilizada por todos proyectos de datos citados anteriormente. Podría decirse que la principal razón por la que la codificación humana siga siendo la práctica habitual es que, en los últimos años, los estudiosos de los conflictos han tratado de construir conjuntos de datos sobre la base de una información cada vez más compleja extraída de los informes de los medios de comunicación (por ejemplo, Hammond y Weidmann 2014). Dicho esto, la codificación automatizada tiene importantes ventajas en comparación con la codificación humana y, por tanto, es probable que adquiera mayor relevancia en los estudios de la conflictividad en los próximos años.

Como entorno analítico, la TEE ofrece una importante base metodológica para los estudios sobre los eventos conflictivos y orienta a los desarrolladores y usuarios de conjuntos de datos antiguos y nuevos. Para los desarrolladores, el marco TEE establece sistemáticamente los diferentes tipos de errores que deben reflejarse en los libros de códigos de datos o artículos que introducen nuevos conjuntos de datos, y apoya la estandarización de las prácticas de información en el campo. De hecho, como la recopilación de datos de eventos y el uso de los datos de los medios de comunicación han sido adoptados por diferentes subdisciplinas y áreas de las ciencias sociales, los tipos de errores de los que se ocupan los investigadores, o de los que informan, parecen

Demarest, Leila y Langer, Arnim "Cómo se incorporan (o no) los eventos en los sets de datos: Los obstáculos y las orientaciones en la utilización de los periódicos en el estudio de los conflictos" *Revista de Estudios Marítimos y Sociales*, N°24, ene-jun 2024, pp. 197-244.





ser muy diferentes. Como también se desprenderá de nuestro análisis del estado de la cuestión, se sabe relativamente poco sobre los errores que pueden surgir al recopilar datos sobre eventos conflictivos en el mundo en desarrollo. Para llenar este vacío, es necesario realizar nuevas investigaciones empíricas. Sobre la base del entorno TEE, pudimos identificar una serie de vías importantes para investigaciones futuras. El entorno TEE también es muy útil para los usuarios de datos sobre conflictos, ya que proporciona información importante sobre la gama de errores que hay que tener en cuenta al utilizar un conjunto de datos específico sobre conflictos, y cómo estos errores pueden afectar a los resultados de la investigación.

En la siguiente sección, desarrollamos el entorno TEE y discutimos en profundidad los errores de medición y representación que pueden surgir durante cada etapa del proceso de investigación. Nuestro análisis se apoya en ejemplos empíricos (necesariamente eclécticos) extraídos de la literatura. En la tercera sección, basada en el entorno TEE, presentamos directrices y estrategias para la recogida de datos de eventos y para futuras investigaciones. En la cuarta sección se presentan las conclusiones.

### **El entorno TEE**

El entorno TEE se inspira en el conocido marco TSE utilizado en la investigación de encuestas (Groves et al. 2004:41-63). En el marco TSE, los errores de medición se producen cuando el valor medido se desvía del valor real. Esto puede surgir de una redacción poco clara de las preguntas y de las escalas de respuesta, de la presencia de un entrevistador, que inhibe al encuestado de responder con sinceridad (es decir, el sesgo de deseabilidad social), o del procesamiento incorrecto de los datos. Los errores de representación se producen cuando no se toman muestras de todas las observaciones existentes en la encuesta. Una característica esencial de una encuesta es el muestreo de sólo una subsección de la población, lo que implica que esta forma de error siempre se produce. Lo importante es que las observaciones se muestreen de forma aleatoria. Esta aleatoriedad puede verse comprometida por un marco de muestreo defectuoso, la falta de respuesta o los ajustes de datos basados en una fuente externa defectuosa (por ejemplo, un censo obsoleto). Evidentemente, pueden producirse dos formas de error: el sesgo, que provoca una desviación sistemática del valor real, y la falta de fiabilidad, que surge de errores aleatorios, lo que hace que los resultados sean menos precisos.



La recopilación y el uso de los datos de los eventos se asemejan al proceso de la encuesta en importantes aspectos. Una similitud importante es que el muestreo es inherente al proceso. Al seleccionar las fuentes de los medios de comunicación para captar los acontecimientos del conflicto, se es consciente de que no se informa necesariamente de todos los acontecimientos que han tenido lugar. El reto surge de los procesos no aleatorios que dirigen la inclusión de eventos en los medios de comunicación, un debate que puede relacionarse con la preocupación por el error de falta de respuesta en las encuestas.

Al igual que los encuestados, las fuentes de noticias y los reportajes pueden presentar la información de forma sesgada o simplemente no pueden proporcionar la información necesaria, lo que da lugar a la falta de datos. Mientras que el entrevistador suele desempeñar un papel clave en el muestreo de los encuestados para las encuestas de opinión pública, lo mismo ocurre con el codificador encargado de muestrear los acontecimientos relevantes en un conjunto de datos. Además, la falta de claridad en las instrucciones de codificación o en las definiciones y categorías de las variables puede dar lugar a datos poco fiables o sesgados, al igual que la falta de claridad en las preguntas de la encuesta. Para ambos tipos de datos, los investigadores pueden intentar validar los datos con una fuente externa. Puede tratarse de un censo o de registros médicos en el caso de las encuestas, o de informes de la policía y de organizaciones no gubernamentales (ONG) en el caso de los datos sobre acontecimientos. Por último, en la fase de análisis, los investigadores pueden optar por ponderar los datos para compensar la falta de respuesta o la selección sesgada o pueden optar por imputar los datos que faltan para preservar el número de casos en el análisis.

La figura 1 muestra el entorno TEE. En el centro de la figura están los pasos de la investigación en la recopilación y el análisis de datos de eventos. Estos pasos no se dan necesariamente de forma secuencial y pueden interactuar de formas importantes. Por ejemplo, el desarrollo del libro de códigos no se finaliza necesariamente antes de codificación, ya que una prueba piloto de codificación ayuda a perfeccionar el libro de códigos. En el caso de la codificación automatizada, el codificador no desempeña ningún papel, o al menos es uno mucho más limitado. El desarrollo del libro de códigos (o diccionario en las aplicaciones automatizadas) adquiere mayor importancia. Además, las comparaciones con fuentes no mediáticas no suelen realizarse, simplemente por la

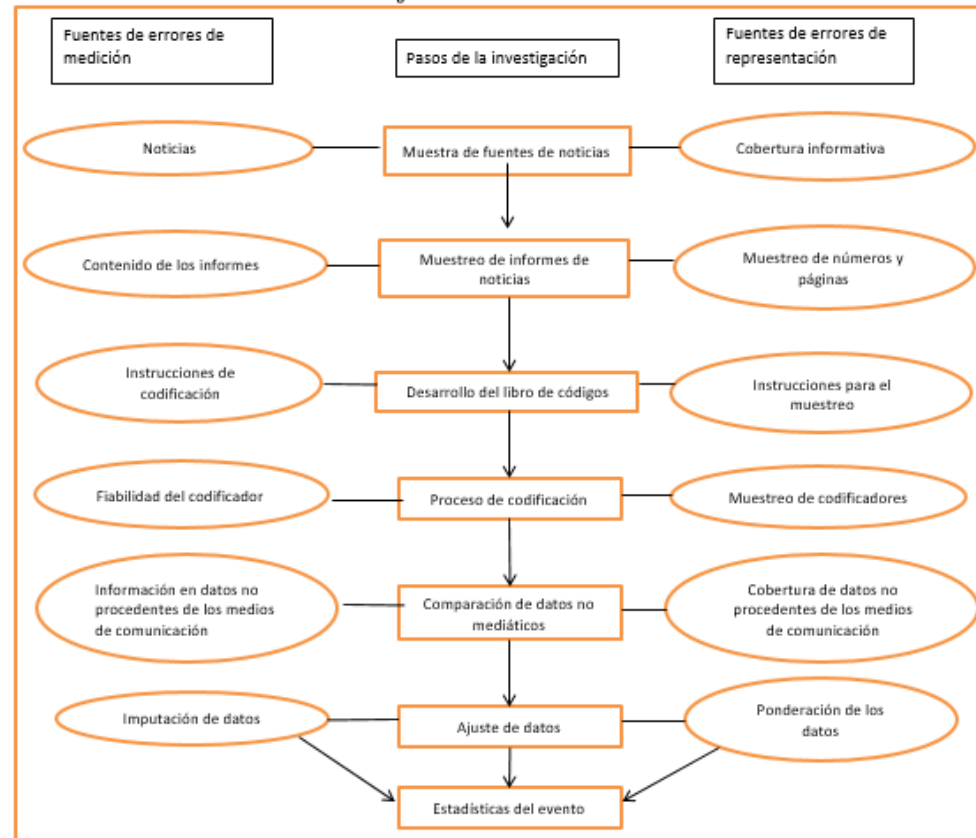
Demarest, Leila y Langer, Arnim "Cómo se incorporan (o no) los eventos en los sets de datos: Los obstáculos y las orientaciones en la utilización de los periódicos en el estudio de los conflictos" *Revista de Estudios Marítimos y Sociales*, N°24, ene-jun 2024, pp. 197-244.





falta de esos datos externos. En línea con el estudio de Groves et al. (2004:48), asociamos cada paso de la investigación con un error de medición y de representación. Ya se han mencionado varias fuentes de error, pero se abordan con mayor profundidad en las secciones siguientes. Estructuramos la discusión según los pasos de investigación identificados.

Figura 1 Errores totales en eventos



Fuente: Demarest y Langer, 2022:637



## Muestreo de fuentes de noticias

Cobertura informativa. El muestreo de fuentes de noticias puede dar lugar tanto a errores de medición como de representación. Comenzamos el debate con el error de representación causado por los efectos de la cobertura de las noticias, ya que se relaciona con el problema relativamente conocido del sesgo de selección (Earl et al. 2004:68-72, Jenkins y Maher 2016, Ortiz et al. 2005). No obstante, aunque el sesgo ha sido ampliamente estudiado en la literatura (Earl et al. 2004:68-72, Jenkins y Maher 2016, Ortiz et al. 2005), los efectos de cobertura también pueden estar asociados a la falta de fiabilidad, al igual que ocurre con el error de muestreo.

A la hora de decidir sobre la recopilación de datos de eventos para tipos específicos de conflicto, períodos de tiempo y entornos geográficos, los investigadores primero deciden sobre los medios de comunicación de la que se extraen los datos. Puede ser un periódico, un servicio de noticias, o incluso las noticias de la televisión y la radio. La elección de una fuente de noticias implica que los eventos incluidos en el conjunto de datos dependen de la selección de los medios de comunicación (o muestreo) de los eventos en las noticias. Como han demostrado multitud de estudios, esta selección dista mucho de ser aleatoria. Se evidencia un doble problema: la cobertura de las fuentes de noticias puede estar determinada por las características de un evento, pero también por las características de la propia fuente de noticias. El primero se considera un efecto de cobertura común a los diferentes medios de comunicación, pero el segundo puede ser específico de la fuente y subraya la importancia de la pregunta de desde qué fuentes de noticias extraer los datos. En primer lugar, analizamos los efectos de selección generales y, a continuación, los específicos de las fuentes.

El trabajo seminal de Galtung y Ruge (1965) sobre la presentación de las crisis del Congo, Cuba y Chipre en los periódicos noruegos sentó las bases de la teoría del valor de las noticias en la ciencia de la comunicación, que investiga las características de un evento que probablemente lo conviertan en noticia (Harcup y O'Neill 2001). Galtung y Ruge proponen 12 factores noticiosos que determinan la difusión de un evento de crisis en el extranjero, como la amplitud o la importancia del evento y la participación de actores de élite. Siguiendo su trabajo, otros científicos de la comunicación han investigado los valores noticiosos que determinan la selección en los medios de

Demarest, Leila y Langer, Arnim "Cómo se incorporan (o no) los eventos en los sets de datos: Los obstáculos y las orientaciones en la utilización de los periódicos en el estudio de los conflictos" *Revista de Estudios Marítimos y Sociales*, N°24, ene-jun 2024, pp. 197-244.





comunicación y han aumentado o reducido el número de factores relevantes (por ejemplo, Harcup y O'Neill 2001, 2016).

Los estudiosos de los movimientos sociales se centran específicamente en los eventos de protesta (Hutter 2014, Koopmans y Rucht 2002). Resumiendo las conclusiones de los estudios anteriores, Earl et al. (2004:69), Ortiz et al. (2005:398-400), Jenkins y Maher (2016:45-46) y Hutter (2014:350-51) señalan que los eventos de protesta a gran escala con muchos participantes, eventos caracterizados por la violencia (daños materiales o físicos, represión policial, detenciones, etc.), eventos organizados por movimientos con personal profesional (de relaciones públicas), y eventos en los que participan actores de alto perfil, son más propensos a ser reportados. Estos hallazgos se apoyan en las comparaciones de la inclusión de eventos entre las fuentes de los medios de comunicación, pero también en las comparaciones de los informes de los medios de comunicación con fuentes externas, como los registros policiales.

El error de representación se ha investigado mucho menos con respecto a los datos de eventos en el mundo en desarrollo. El reciente interés por los conflictos de bajo nivel incluyendo las protestas en los países en vías de desarrollo, quizás pueda asumir las mismas preferencias de cobertura. En el caso de los eventos de conflictos armados, podríamos suponer que debido al nivel de violencia, su incorporación en las noticias es muy probable. Sin embargo, los contextos en los que surgen los conflictos armados suelen ser diferentes de los escenarios occidentales que se suelen investigar. Las guerras civiles suelen estallar en zonas rurales alejadas del centro de poder del gobierno (por ejemplo, Kalyvas 2006:38-48), lo que repercute en la infraestructura de comunicaciones presente en la región. Además, aunque los conflictos armados atraen la atención periodística, el clima de violencia y los daños en las infraestructuras pueden obstaculizar la cobertura de los eventos. A este respecto, un estudio de Weidmann (2016) es muy instructivo. Compara un conjunto de datos sobre conflictos armados en Afganistán recogidos por el ejército estadounidense -y revelados por WikiLeaks- con el conjunto de datos de la GED de la UCDP (que sólo utilizó fuentes de los medios de comunicación para este conflicto) y afirma que la cobertura de los teléfonos móviles aumenta significativamente la probabilidad de que los medios de comunicación

Demarest, Leila y Langer, Arnim "Cómo se incorporan (o no) los eventos en los sets de datos: Los obstáculos y las orientaciones en la utilización de los periódicos en el estudio de los conflictos" *Revista de Estudios Marítimos y Sociales*, N°24, ene-jun 2024, pp. 197-244.





informen de los eventos, lo que sugiere una subrepresentación sistemática de los eventos en zonas rurales remotas.

Los efectos de la selección específica de la fuente están relacionados con la orientación ideológica u orientación política de una fuente de noticias, así como su alcance geográfico (p. ej, Davenport 2010:107-26). Por ejemplo, en el análisis de eventos de protesta varios estudios constatan que los periódicos conservadores reportan de manera insatisfactoria las manifestaciones violentas para limitar el comportamiento de los imitadores (para una visión general, véase Ortiz et al. 2005:401). El segundo factor, el ámbito geográfico de la fuente de noticias se refiere a si se llega a un público local, nacional o internacional. Dedicamos aquí más atención a esta cuestión, ya que muchos conjuntos recientes de datos hacen uso de inventarios de múltiples fuentes, como Factiva, LexisNexis o Keesing's Record of World Events, que se basan en gran medida en los servicios de noticias internacionales para codificar los eventos que se producen en una amplia gama de países en desarrollo, incluidos los conflictos armados violentos y las protestas.

Las fuentes de noticias locales pueden cubrir los conflictos locales más ampliamente que las fuentes nacionales, que aplican un procedimiento de selección adicional. Las fuentes internacionales tienen un proceso de selección aún más estricto. Sin embargo, para algunos eventos, como los conflictos armados en curso, los servicios profesionales de noticias internacionales podrían ser potencialmente más valiosos que los servicios de los medios de comunicación locales (interrumpidos). La forma en que el sesgo de selección se produce cuando las escalas de conflicto y el alcance de las fuentes de noticias interactúan es una cuestión empírica muy relevante y quizás insuficientemente tratada. Varios estudios indican su importancia. Herkenrath y Knoll (2011), por ejemplo, encuentran que los periódicos internacionales informan sustancialmente menos sobre eventos de protesta que los periódicos nacionales en Argentina, México y Paraguay y que estas diferencias están relacionadas con, entre otras cosas, el uso de la violencia, pero también con una diferencia general en la atención de los medios internacionales hacia estos tres países.

Bueno de Mesquita et al. (2015) desarrolló un conjunto de datos sobre la violencia política en Pakistán basado en los periódicos nacionales y registran un mayor número de

Demarest, Leila y Langer, Arnim "Cómo se incorporan (o no) los eventos en los sets de datos: Los obstáculos y las orientaciones en la utilización de los periódicos en el estudio de los conflictos" *Revista de Estudios Marítimos y Sociales*, N°24, ene-jun 2024, pp. 197-244.





incidentes que los conjuntos de datos que se basan en Factiva. Demarest y Langer (2018) compararon conjuntos de datos basados en fuentes de noticias internacionales frente a nacionales sobre eventos de conflicto en Nigeria. Encuentran que las fuentes internacionales subrepresentan los eventos de conflicto, en particular los eventos de protesta. En ambos estudios también se observa que el subregistro relativo afecta a la distribución subnacional de los eventos, una línea de investigación cada vez más importante en los estudios sobre conflictos (por ejemplo, Gleditsch et al. 2014:303-05). Por último, Barron y Sharpe (2008) utilizaron fuentes de noticias a nivel de distrito para captar eventos violentos en Indonesia y muestran que estas registran más incidentes que los periódicos provinciales y, por tanto, proporcionan una mayor comprensión de las causas locales del conflicto.

En general, se argumenta que los inventarios de múltiples fuentes son más fiables que fuentes individuales (Jenkins y Maher 2016:47-49), pero es importante tener en cuenta que los inventarios multifuente no incluyen el "universo de los medios de comunicación" (Ortiz et al. 2005:402). Especialmente en el caso de los conflictos en el mundo en desarrollo, es importante tener en cuenta que las fuentes internacionales en inglés incluidas en estos inventarios podrían no cubrir suficientemente estos escenarios (Schrodt 2012:552-53). En principio, los procedimientos de codificación automatizada no son sensibles a los errores de representación, pero requieren un texto legible por máquina y se inspiran predominantemente en servicios internacionales de información de noticias como Reuters y LexisNexis (por ejemplo, el conjunto de datos Integrated Data for Events Analysis, Bond et al. 2003, el conjunto de datos Global Data on Events, Location and Tone [GDELT], Leetaru y Schrodt 2013, el conjunto de datos Kansas Event Data System [KEDS], Schrodt 2006), lo cual es una característica importante a tener en cuenta. El uso de los periódicos locales para investigar eventos conflictivos en los países en desarrollo es probable que surja como una importante línea de investigación en los estudios sobre conflictos, entre otras cosas porque los periódicos locales están cada vez más disponibles en línea (por ejemplo, el repositorio AllAfrica). Sin embargo, el uso de fuentes locales conlleva nuevos retos, por ejemplo, relacionados con quién es el dueño de los medios de comunicación y el control estatal de los mismos.

Demarest, Leila y Langer, Arnim "Cómo se incorporan (o no) los eventos en los sets de datos: Los obstáculos y las orientaciones en la utilización de los periódicos en el estudio de los conflictos" *Revista de Estudios Marítimos y Sociales*, N°24, ene-jun 2024, pp. 197-244.





Hasta ahora, parece que se han realizado pocas investigaciones sistemáticas para evaluar el impacto de estos problemas en la calidad de los datos sobre eventos conflictivos.

### **La cobertura de las noticias.**

Pasamos ahora al problema del error de medición derivado de la fuente de noticias, que refiere a la información que las fuentes de noticias comunican sobre un evento. Esta forma de error puede relacionarse con el concepto de sesgo de descripción (Earl et al. 2004:72-73). Al hablar de los problemas de sesgo en la descripción varios estudiosos distinguen entre "noticias duras" y "noticias blandas" y sostienen que las primeras están menos sujetas a sesgos que las segundas (Earl et al. 2004:72, Franzosi 1987:7, Raleigh et al. 2010:656). Se sugiere que las noticias duras incluyen el "quién, qué, cuándo, dónde y por qué del evento" (Earl et al. 2004:72), mientras que las noticias blandas incluyen las interpretaciones de las causas y las consecuencias, las representaciones de los actores, etc. En primer lugar, discutimos las investigaciones sobre las dimensiones de las noticias blandas y, luego, nos centramos en las noticias duras. Sin embargo, como se argumentará más adelante, la distinción que se hace en la literatura entre noticias duras y blandas no es sencilla. Además, no todas las inexactitudes de las noticias son necesariamente signos de parcialidad, sino que también pueden indicar falta de fiabilidad debido a las dificultades que tienen las fuentes de los medios de comunicación para adquirir determinados tipos de información.

Los efectos de las noticias blandas pueden relacionarse con el concepto de encuadre (framing en inglés). Existen varias definiciones de encuadre, como ejemplo, citamos a Entman (1993:52). "Enmarcar es seleccionar algunos aspectos de una realidad percibida y hacerlos más destacados en un contexto de comunicación, de tal manera que se promueva una determinada definición del problema, una interpretación causal, una evaluación moral y/o una recomendación de tratamiento". Una gran cantidad de investigaciones se han centrado en la forma en que los medios de comunicación representan las acciones de protesta, aunque se centran predominantemente en entornos occidentales (por ejemplo, Dardis 2006, McLeod y Hertog 1992). Una importante línea de investigación se centra en las diferencias de encuadre según el perfil ideológico



(conservador o liberal) de la fuente de noticias y en si los periódicos conservadores son más propensos a representar negativamente a los manifestantes (por ejemplo, Chan y Lee 1984, Lee 2014, Weaver y Scacco 2012).

Aunque la literatura sobre el encuadre de los eventos conflictivos ofrece interesantes ideas sobre las orientaciones de las diferentes fuentes de noticias, no está tan claro hasta qué punto las diferentes representaciones del conflicto pueden afectar a los conjuntos de datos de eventos que se centran en fechas, lugares y actores. No obstante, la línea que separa las noticias blandas y las noticias duras no es necesariamente clara. Un ejemplo bien conocido de "dato duro" comúnmente disputado es el número de participantes en una protesta, que puede ser exagerado por los activistas o subestimado por las autoridades policiales (Day, Pinckney y Chenoweth 2015:130). Además, las estimaciones de víctimas mortales se consideran a menudo como hechos concretos que son difíciles de establecer (Raleigh et al. 2010:656, Sundberg y Melander 2013:527). La fuente -policía frente a manifestantes o gobierno frente a rebeldes- que prefiere un periódico de la muestra puede sesgar las estadísticas de los datos de los eventos.

Relativamente pocos estudios han comparado los datos externos con la información de hechos concretos en los medios de comunicación. McCarthy et al. (1999:117-26) compararon datos de los registros policiales con los informes de los medios de comunicación impresos (y electrónicos) sobre eventos de protesta en Washington, DC. Encontraron una buena correspondencia en cuanto a las fechas de las protestas y de los objetivos, pero una menor coincidencia en cuanto al tamaño de las mismas. Sin embargo, esto último podría estar relacionado con el sesgo o la falta de fiabilidad, ya que el análisis no describe cómo se relacionan las variables entre sí.

Weidmann (2015) investigó las diferencias en la notificación de hechos concretos entre el conjunto de datos militares de EE.UU. sobre eventos de conflictos armados en Afganistán y el UCDP GED. Afirma que, para la mayoría de los eventos, el número de víctimas del conjunto de datos militares se encuentra dentro de la estimación de bajas del conjunto de datos del UCDP GED. Sin embargo, hay más eventos para los que el UCDP GED da una estimación más alta, lo que podría indicar un ligero sesgo hacia la presentación de cifras de víctimas más altas en las noticias. También hay diferencias entre el UCDP GED y el conjunto de datos militares en cuanto a la localización del

Demarest, Leila y Langer, Arnim "Cómo se incorporan (o no) los eventos en los sets de datos: Los obstáculos y las orientaciones en la utilización de los periódicos en el estudio de los conflictos" *Revista de Estudios Marítimos y Sociales*, N°24, ene-jun 2024, pp. 197-244.







evento. Basándose en sus análisis, Weidmann (2015:1143) sostiene que los investigadores no deberían utilizar los datos para los análisis por debajo de un rango de 50 km. Su investigación indica que incluso un dato duro como la "ubicación" tampoco se comunica siempre de forma fiable, en particular cuando se trata de eventos de conflictos armados. Aunque las víctimas mortales se consideran difíciles de establecer de forma fiable, la investigación de Weidmann sugiere que su registro parece estar relativamente libre de errores.

La cobertura informativa y los informes están relacionados con algunos de los errores más conocidos descritos en la literatura. No obstante, muchos estudios se centran en contextos occidentales y eventos de protesta, y sólo en menor medida en contextos en desarrollo y eventos de conflicto violento (armado). Aunque se pueden extraer importantes principios y lecciones de los movimientos sociales y los estudios de comunicación, es evidente la necesidad de más investigación empírica sobre estas formas de errores en los estudios de conflictos. En las siguientes secciones, nos centraremos en los errores que posiblemente se discutan menos en los estudios actuales. Estos errores no surgen necesariamente del funcionamiento de los medios de comunicación, sino que están más relacionados con los procedimientos de recogida de datos.

### **Muestreo de informes de noticias**

Muestreo de números y páginas. Mientras que algunos investigadores recurren a todos los informes disponibles de una fuente específica, otros se basan en el muestreo adicional de números o páginas de periódicos específicos (Earl et al. 2004:68, Krippendorff 2013:112-25). Cuando se trata de conjuntos de datos de conflictos recientes (véase la Introducción), no se incluye esta etapa de muestreo adicional, ya que normalmente se basan en informes extraídos de inventarios de múltiples fuentes, utilizando términos clave y especificaciones de fechas y países. En el caso de los estudios que se basan en periódicos nacionales o locales, especialmente cuando se estudia un periodo relativamente extenso, este muestreo adicional puede ser necesario para reducir los costes de codificación. Por ejemplo, cuando Kriesi et al. (1998:253-63) realizaron un estudio fundamental sobre eventos de protesta en cuatro países de Europa Occidental, utilizaron un periódico nacional por país, pero sólo la edición del lunes.

Demarest, Leila y Langer, Arnim "Cómo se incorporan (o no) los eventos en los sets de datos: Los obstáculos y las orientaciones en la utilización de los periódicos en el estudio de los conflictos" *Revista de Estudios Marítimos y Sociales*, N°24, ene-jun 2024, pp. 197-244.





Cubrieron el periodo comprendido entre 1975 y 1989. Incluso si no hay sesgos sistemáticos asociados a las ediciones específicas de los periódicos, este muestreo adicional genera una mayor falta de fiabilidad. También es posible seleccionar sólo la primera página de un número de periódico, lo que podría, por ejemplo, reforzar el sesgo hacia la inclusión de eventos de alto perfil caracterizados por la violencia.

### **El contenido del informe.**

Las preferencias periodísticas o editoriales también pueden ser una importante fuente de error a nivel de la noticia. Por ejemplo, Chojnacki et al. (2012:390-92) llaman la atención sobre el hecho de que los informes de noticias suelen citar fuentes que tienen incentivos para proporcionar información sesgada. Utilizan el ejemplo de un informe en el que un líder rebelde afirmaba haber matado a 30 soldados del gobierno. En su conjunto de datos (Event Data on Conflict and Security [EDACS]), crearon una variable adicional, indicando que la información podría estar sesgada si se utilizan fuentes dudosas.

Además de los sesgos que pueden surgir de las preferencias informativas, las noticias pueden ser fuentes importantes de falta de fiabilidad. En primer lugar, los informes sobre los eventos pueden ser detallados o vagos. Algunos informes pueden proporcionar información sobre el tamaño de un grupo de manifestantes, mientras que otro informe sobre el mismo evento puede sólo mencionar que se produjo la protesta. Del mismo modo, se puede informar de la captura de un territorio por un grupo rebelde, pero no necesariamente de si hubo víctimas. En algunos casos, múltiples informes pueden proporcionar información adicional valiosa, pero en otros, los informes imprecisos pueden ser la única fuente de información y el resultado puede ser un grado considerable de falta de datos. El valor periodístico de un evento también puede afectar a la profundidad de la información y a la extensión del artículo de prensa dedicado al mismo. Por ejemplo, una protesta a gran escala puede atraer más atención informativa que un evento de tamaño limitado, y por lo tanto, también puede aparecer más información sobre el evento. No obstante, aunque algunos eventos pueden atraer una gran atención informativa, como las graves violaciones de los derechos humanos en los conflictos armados, la "niebla de la guerra" también puede impedir la obtención de información fiable.

Demarest, Leila y Langer, Arnim "Cómo se incorporan (o no) los eventos en los sets de datos: Los obstáculos y las orientaciones en la utilización de los periódicos en el estudio de los conflictos" *Revista de Estudios Marítimos y Sociales*, N°24, ene-jun 2024, pp. 197-244.





En segundo lugar, los informes también pueden arrojar dudas explícitas sobre si un evento ocurrió y cómo, sobre la identidad de los actores o sobre la validez de una estimación de víctimas. Los informes pueden, por ejemplo, afirmar que la identidad de los atacantes o presuntos rebeldes es incierta. Estas formas de error de medición sólo se recogen si se incluyen estas variables indicadoras en el libro de códigos.

En tercer lugar, los problemas de codificación pueden surgir de informes contradictorios. Aunque el carácter incompleto de los informes de noticias lleva a muchos investigadores a basarse en múltiples informes para construir variables de datos de eventos, esto también puede plantear preguntas adicionales sobre la forma en que se combinan los informes (Weidmann y Rød 2015:125-26). Un problema crucial surge cuando la información es inconsistente. Algunos conjuntos de datos proporcionan instrucciones a los codificadores para que agreguen la información de formas particulares. Por ejemplo, el SCAD establece que en el caso de estimaciones de víctimas múltiples, se toma la media (Codebook version 3.1.), mientras que el ACLED establece que se debe utilizar el número más bajo (Raleigh y Dowd 2017:20). Otras soluciones a los informes contradictorios sugieren codificar cada informe individualmente. Basándose en su trabajo sobre los eventos de protesta para el conjunto de datos de Campañas y Resultados No Violentos y Violentos, Day et al. (2015:130-31) recomiendan la codificación de diferentes informes, junto con la inclusión de una variable de rango de ambigüedad métrica en el conjunto de datos finales del evento.

Del mismo modo, Weidmann y Rød (2015) proponen la creación de un conjunto de datos intermedio, que incluye la codificación de eventos por informe de noticias, y un conjunto de datos de eventos, que agrega la información a través de los informes. Como se proporciona toda la información de los informes, las reglas de agregación (media, mínimo, etc.) pueden modificarse. La codificación de los informes de noticias por separado puede aumentar la transparencia y la posibilidad de réplica, a diferencia de permitir que los codificadores agreguen ellos mismos los informes de noticias. Esta elección de codificación también puede tener importantes implicaciones para el control del proceso de codificación y las puntuaciones de fiabilidad entre codificadores (véase más adelante). Sin embargo, codificar las noticias por separado puede aumentar los costes de la investigación.

Demarest, Leila y Langer, Arnim "Cómo se incorporan (o no) los eventos en los sets de datos: Los obstáculos y las orientaciones en la utilización de los periódicos en el estudio de los conflictos" *Revista de Estudios Marítimos y Sociales*, N°24, ene-jun 2024, pp. 197-244.





## Desarrollo del libro de códigos

### Instrucciones de muestreo.

Las instrucciones del libro de códigos son cruciales para evitar la confusión de los codificadores y para apoyar un muestreo coherente, así como la codificación de los eventos relevantes. Cuando se utiliza la codificación automática, el diccionario y el programa de codificación determinan la selección y codificación de los casos en el conjunto de datos basándose en la identificación de los actores relevantes, las palabras, etc., en lugar de un codificador.<sup>6</sup> Generalmente, los libros de códigos y los diccionarios se revisan después de una fase inicial de prueba de codificación, en la que se revelan las posibles fuentes de error. Las instrucciones de muestreo son una preocupación importante: ¿Qué eventos deben incluirse en el conjunto de datos y cuáles deben excluirse?

Al elaborar las instrucciones para los codificadores humanos, los investigadores pueden adoptar una definición o proporcionar una lista de eventos elegibles (por ejemplo, Kriesi et al. 1998:263-69). Muchos conjuntos de datos de eventos conflictivos se basan principalmente en definiciones de eventos, pero una advertencia potencial es que cuanto más estricta sea la definición, más difícil será codificar de forma coherente los informes de eventos imprecisos. Por ejemplo, los informes no siempre dan detalles sobre los actores, qué actor utilizó la violencia o el número de participantes en un evento, lo que puede crear confusión e incoherencias de muestreo cuando la categorización requiere esta información. También puede ser importante incluir instrucciones sobre cómo tratar los casos en los que un informe arroja dudas sobre su ocurrencia o elegibilidad para su inclusión.

Cuando se toman muestras de eventos de repositorios en línea, se aplican las mismas preocupaciones. En bases de datos como LexisNexis, se puede desarrollar una cadena de búsqueda de palabras clave relevantes y aplicarlas para extraer noticias sobre un tema o evento. Después, los codificadores pueden verificar manualmente una submuestra para seleccionar la utilidad y eficacia de la cadena de búsqueda y la cantidad de "ruido". No obstante, la decisión del codificador de incluir o excluir eventos sigue

Demarest, Leila y Langer, Arnim "Cómo se incorporan (o no) los eventos en los sets de datos: Los obstáculos y las orientaciones en la utilización de los periódicos en el estudio de los conflictos" *Revista de Estudios Marítimos y Sociales*, N°24, ene-jun 2024, pp. 197-244.



requiriendo consistencia y replicabilidad y la consideración de las cuestiones antes mencionadas. Una noticia que incluya palabras clave relevantes como "violencia" puede informar de más de un evento, por ejemplo, todos los cuales deben ser muestreados de forma coherente. Además, el uso de cadenas de búsqueda no garantiza que todos los eventos muestreados de una fuente de noticias también lo sean utilizando términos específicos. Aunque las cadenas de búsqueda suelen incluir muchos términos clave, algunos eventos pueden pasar desapercibidos.

Los problemas de ruido y la omisión de eventos son también una preocupación importante cuando se utilizan procedimientos de codificación automatizados. Sin embargo, es útil señalar primero las ventajas de la codificación automática. Aunque el desarrollo de diccionarios requiere tiempo, incluyendo el mayor número posible de verbos y frases clave, variaciones, nombres de actores (por ejemplo, Estados Unidos, EE.UU., presidente Trump), una vez desarrollados, ofrecen la posibilidad de recorrer grandes volúmenes de datos en segundos (Bond et al. 1997, Schrodt y Van Brackle 2013). Además, una revisión del diccionario no conlleva un proceso de recodificación que requiere mucho tiempo. En su lugar, el programa puede volver a ejecutarse sobre los mismos datos con el diccionario revisado. Por último, los diccionarios pueden compartirse entre investigadores y utilizarse para nuevos proyectos. Sin embargo, un punto importante de debate es si la codificación automática es capaz de identificar los eventos "correctos" y si estos eventos se codifican correctamente (véase más adelante), ya que la codificación humana suele tomarse como norma.

La capacidad de los procedimientos de codificación automática para incluir un número suficientemente alto de eventos relevantes ("recuperación"), excluyendo al mismo tiempo los eventos irrelevantes (" precisión ") -los eventos relacionados con competiciones deportivas son falsos positivos- es un importante reto de muestreo.<sup>7</sup> Varios investigadores han investigado empíricamente la recuperación y/o la precisión de las aplicaciones de codificación automática en comparación con un set de entrenamiento desarrollado por codificadores humanos. Por ejemplo, Bond et al. (1997) afirman que el programa original de análisis sintáctico disperso de KEDS 8 funciona al menos tan bien como los (nuevos) codificadores humanos en la identificación de eventos relevantes (alrededor del 80%). King y Lowe (2003) prueban el lector VRA y

Demarest, Leila y Langer, Arnim "Cómo se incorporan (o no) los eventos en los sets de datos: Los obstáculos y las orientaciones en la utilización de los periódicos en el estudio de los conflictos" *Revista de Estudios Marítimos y Sociales*, N°24, ene-jun 2024, pp. 197-244.





afirman que su rendimiento es tan bueno como el de los codificadores humanos en lo que respecta a la recuperación (93%), pero menos en cuanto a la precisión (23% de aciertos). Sin embargo, en general, son positivos respecto al potencial de la codificación automática.

Además de las comparaciones con codificadores humanos, también se han realizado comparaciones entre programas que están en continuo desarrollo. Boschee, Natarajan y Weischedel (2013) comparan el programa TABARI desarrollado por Schrodts como continuación del programa KEDS original y afirman que, en lo que respecta a la recuperación y la precisión, su procedimiento de análisis sintáctico disperso es significativamente superado por el programa SERIF de BBN que se basa en el procesamiento del lenguaje natural. 9 Más recientemente, Croicu y Weidmann (2015) desarrollaron un sistema clasificador de aprendizaje automático que muestra porcentajes de recuperación y precisión de alrededor de 90 y 50, respectivamente, de nuevo en comparación con los codificadores humanos. Heap et al. (2017) proponen un proceso conjunto humano/máquina para la selección de texto relevante mediante aprendizaje automático supervisado para mejorar la recuperación y precisión. Además del procesamiento del lenguaje natural y el aprendizaje automático, otra área de progreso en la codificación automatizada reside en los campos aleatorios condicionales (Schrodts y Van Brackle 2013:38, Stepinski, Stoll y Subramanian 2006).

Parece que la codificación automatizada tiene importantes y crecientes beneficios para el muestreo de eventos. Sin embargo, sigue habiendo una serie de retos que hay que tener en cuenta. El primer reto crucial se refiere a la duplicación o a la inclusión del mismo evento varias veces en el conjunto de datos (Bond et al. 2003:737-38, Schrodts y Van Brackle 2013). Todavía no existe un procedimiento automático real para filtrar los duplicados, excepto para descartar eventos con la misma hora, ubicación, actores, etc. La revisión humana del conjunto de datos puede ser necesaria para excluir más duplicados y puede ser un ejercicio costoso cuando se trata de grandes volúmenes de datos. Otro reto tiene que ver con el idioma, ya que la mayoría de los diccionarios y aplicaciones se centran predominantemente en el idioma inglés (Schrodts y Van Brackle 2013:45), mientras que las extensiones a otros idiomas pueden llevar a la inclusión de fuentes más diversas y no occidentales. Sin embargo, el uso del inglés tampoco es

Demarest, Leila y Langer, Arnim "Cómo se incorporan (o no) los eventos en los sets de datos: Los obstáculos y las orientaciones en la utilización de los periódicos en el estudio de los conflictos" *Revista de Estudios Marítimos y Sociales*, N°24, ene-jun 2024, pp. 197-244.





uniforme, y la elección de palabras y las estructuras de las frases también pueden variar según la región o el país, y pueden ser más marcadas para los eventos nacionales que para los internacionales (Schrodt, Simpson y Gerner 2001). Incluso la propia fuente de noticias puede variar en el uso del lenguaje (Boschee et al. 2013).

La codificación automatizada se ha utilizado principalmente para la recopilación de eventos políticos en el campo de las relaciones internacionales (por ejemplo, Schrodt 2006). Cada vez más, el uso de la codificación automatizada se utiliza también para investigar los conflictos internos, incluso en contextos de países en desarrollo (por ejemplo, Leetaru y Schrodt 2013). Esto implica que los desafíos con respecto a la construcción de diccionarios descritos anteriormente son cada vez más pertinentes de tratar, tanto en lo que respecta a la selección de eventos como a la codificación de los mismos, como se verá más adelante.

### **Instrucciones de codificación.**

Las instrucciones de codificación poco claras pueden crear errores de representación y de medición. De nuevo, surge el problema de la definición de los eventos. Por ejemplo, el libro de códigos de la USDAA incluye 12 definiciones de eventos, pero se afirma que estos tipos de conflicto "no son en absoluto categorías mutuamente excluyentes [ . . . ] Aunque hemos tratado de ser coherentes en la codificación de dichos eventos, hay que tener cuidado al tratar las categorías como fenómenos claramente distinguibles" (Urdal 2008:11). Este problema se deriva de la falta de información o de la información conflictiva en los informes de los eventos. Por ejemplo, en algunos conjuntos de datos, la mención de una asociación detrás de la protesta puede marcar la diferencia entre la categorización como protesta espontánea o como protesta organizada (por ejemplo, SCAD).<sup>10</sup> Sin embargo, esto también puede estar influenciado por la profundidad de la información.

Al elaborar el libro de códigos, los investigadores pueden tener que elegir entre categorías muy genéricas de eventos, actores, etc., que pueden codificarse de forma fiable, o categorías muy específicas, para las que la codificación es menos fiable. Se trata de un intercambio importante que hay que hacer. Aunque las categorías amplias o genéricas pueden crear una mayor coherencia, puede que no proporcionen el nivel de

Demarest, Leila y Langer, Arnim "Cómo se incorporan (o no) los eventos en los sets de datos: Los obstáculos y las orientaciones en la utilización de los periódicos en el estudio de los conflictos" *Revista de Estudios Marítimos y Sociales*, N°24, ene-jun 2024, pp. 197-244.





precisión de la información que los investigadores buscan. Por ejemplo, una categoría genérica de actores como "atacantes" podría codificarse de forma muy fiable, pero también se querría saber, en la medida de lo posible, si los atacantes eran grupos rebeldes concretos o milicias étnicas, partidos políticos, etc. Desgraciadamente, la necesidad de disponer de información detallada sobre los eventos para responder a determinadas preguntas de investigación no siempre se ve satisfecha por la información proporcionada por los medios de comunicación.

Para la codificación automatizada, la complejidad de la codificación de eventos no sólo se ve dificultada por la información disponible en las noticias, sino también por el diccionario y los matices que pueden captar las estructuras oracionales predefinidas. Bond et al. (1997) también analizaron la categorización de eventos, además del muestreo de eventos, y volvieron a encontrar que la codificación automática tiene un rendimiento similar al de la codificación humana. King y Lowe (2003) obtienen afirmaciones similares, pero también muestran que las clasificaciones de eventos más generales se codifican con mayor fiabilidad que las detalladas. El hecho de que las definiciones detalladas de los eventos no siempre son viables es también discutido por Schrodt y Van Brackle (2013:33).

En general, se considera que la codificación automatizada funciona mejor cuando las variables que hay que extraer no son demasiado complejas. Uno de los retos es que el campo de los estudios sobre la paz y los conflictos está avanzando cada vez más hacia definiciones y características de eventos más complejas, así como hacia la recopilación detallada de información sobre el tiempo y la ubicación. Como ya se ha comentado, en la investigación empírica se busca cada vez más información sobre la ubicación subnacional de los eventos, aunque se argumenta que la codificación automatizada funciona mejor a nivel de país (Bond et al. 2003:739, Schrodt y Van Brackle 2013:46). Hammond y Weidmann (2014), por ejemplo, sostienen que el conjunto de datos del GDELT debe utilizarse con precaución para los análisis subnacionales, ya que difiere sustancialmente de la codificación humana y parece mostrar un sesgo hacia las capitales de los países. Hickler y Wiesel (2012) son más optimistas cuando comparan la información espacial de los datos codificados por humanos y por máquinas en el marco del conjunto de datos EDACS, aunque siguen planteando dudas.

Demarest, Leila y Langer, Arnim "Cómo se incorporan (o no) los eventos en los sets de datos: Los obstáculos y las orientaciones en la utilización de los periódicos en el estudio de los conflictos" *Revista de Estudios Marítimos y Sociales*, N°24, ene-jun 2024, pp. 197-244.







Cuando se utiliza la codificación humana, el desarrollo del libro de códigos es un comienzo importante, pero la forma de aplicarlo es en gran medida responsabilidad de los codificadores. La codificación automática excluye a los codificadores o les otorga un papel más limitado (de supervisión) (por ejemplo, Heap et al. 2017). En la siguiente sección, nos centraremos en los errores derivados del codificador en un proyecto típico de codificación humana. Curiosamente, aunque la codificación automática se suele comparar con una codificación humana de referencia, la propia codificación humana también está sujeta a errores sustanciales. De hecho, éste es el argumento principal de Bond et al. (1997:555), que desde un principio lamentaron la mala calidad de la codificación humana.

### Proceso de codificación

**Muestreo del codificador.** Siguiendo las instrucciones del libro de códigos, los codificadores muestrean los eventos en el conjunto de datos y extraen información sobre las variables clave. Por lo tanto, el codificador también puede ser una fuente de error de representación y de medición, y pueden surgir tanto la falta de fiabilidad como el sesgo. En el muestreo, los codificadores pueden pasar por alto eventos completamente al azar debido, por ejemplo, a la falta de atención. El sesgo se produce cuando los codificadores pasan por alto habitualmente ciertos eventos o malinterpretan regularmente las instrucciones sobre lo que constituye un evento relevante. Por desgracia, es probable que eventos más pequeños y de baja escala pasan más a menudo desapercibidos que los eventos de alto perfil anunciados en los titulares (por ejemplo, Kriesi et al. 1998:270), por lo que el error de muestreo del codificador puede reforzar potencialmente el sesgo de selección. El error de muestreo de los codificadores no suele medirse (ni comunicarse), pero algunos investigadores han intentado cuantificarlo. En su trabajo sobre los movimientos sociales en cuatro países de Europa Occidental, Kriesi et al. (1998:270) informan de que en las comparaciones por parejas, alrededor del 60% de los eventos de protesta fueron registrados por ambos codificadores. En el siguiente proyecto, se alcanzó un acuerdo de identificación de alrededor del 70% entre los codificadores (Hutter 2014:355). Aunque utilizaron una fuente de datos diferente a la de los informes de noticias -informes del Secretario General de las Naciones Unidas sobre

Demarest, Leila y Langer, Arnim "Cómo se incorporan (o no) los eventos en los sets de datos: Los obstáculos y las orientaciones en la utilización de los periódicos en el estudio de los conflictos" *Revista de Estudios Marítimos y Sociales*, N°24, ene-jun 2024, pp. 197-244.





operaciones de mantenimiento de la paz-, Ruggeri, Gizelis y Dorussen (2011:348-51) también observan un grave error de muestreo de los codificadores. Constatan que los codificadores independientes solo identificaron dos veces el 18-41% de los eventos relevantes. Para ello, el equipo de investigación tuvo que cambiar de estrategia y hacer que los jefes de equipo identificaran y destacaran eventos relevantes, que luego fueron codificados por los asistentes.

### **Fiabilidad del codificador.**

Para la investigación que hace uso del análisis de contenido de los medios de comunicación, el de la fiabilidad entre codificadores para indicar el error de medición se considera como un imperativo metodológico en la ciencia de la comunicación (Krippendorff 2013:272-73).<sup>11</sup> Este imperativo también se ha introducido en los análisis de eventos de protesta en los estudios de los movimientos sociales (occidentales) (Hutter 2014:354-55). Sin embargo, muchos conjuntos de datos de eventos de conflicto centrados en el mundo en desarrollo no informan de tales mediciones (Ruggeri et al. 2011:356-59, Salehyan 2015:107- 08). No obstante, al realizar pruebas de fiabilidad entre codificadores, se puede comprobar si el mismo instrumento de medición (el libro de códigos) lleva a codificadores independientes a alcanzar resultados similares (Krippendorff 2013:273-75). Las medidas más comunes son la  $\alpha$  de Krippendorff y la  $k$  de Cohen, que corrigen la concordancia fortuita ponderando más la inconsistencia en las categorías de respuesta menos frecuentes en el coeficiente final. Las comprobaciones de fiabilidad entre codificadores pueden utilizarse para perfeccionar el libro de códigos o seleccionar a los "mejores" codificadores después de una fase piloto. Se recomienda realizar pruebas con regularidad a lo largo del proceso de codificación, ya que sólo la realización de pruebas posteriores a la recogida de datos puede revelar la necesidad de descartar o recodificar una cantidad importante de datos. Las pruebas pueden realizarse sobre un pequeño subconjunto de los datos (5-10 por ciento).

Al interpretar la fiabilidad entre codificadores, también es importante tener en cuenta que la fiabilidad entre codificadores puede ser baja si cada codificador comete errores aleatorios (falta de fiabilidad del codificador) o si cada codificador interpreta habitualmente las reglas de forma diferente (sesgo del codificador). Sin embargo, si

Demarest, Leila y Langer, Arnim "Cómo se incorporan (o no) los eventos en los sets de datos: Los obstáculos y las orientaciones en la utilización de los periódicos en el estudio de los conflictos" *Revista de Estudios Marítimos y Sociales*, Nº24, ene-jun 2024, pp. 197-244.





todos los codificadores interpretan erróneamente una regla de codificación, este sesgo no será captado por la estadística de fiabilidad. El cálculo de las estadísticas de fiabilidad entre codificadores puede ser especialmente importante para la investigación de las interpretaciones causales o el encuadre en los informes de los medios de comunicación. No obstante, no es necesariamente seguro suponer que los hechos concretos se codifican relativamente libres de errores (por ejemplo, Eck 2012:130-35).

Por último, cabe señalar que la fiabilidad entre codificadores suele calcularse a nivel de la noticia en los estudios de comunicación. En efecto, este nivel permite un seguimiento más estrecho del trabajo de los codificadores. Sin embargo, cuando los codificadores tienen instrucciones de agregar informes e información de eventos, este proceso de control puede complicarse. Pueden surgir retos importantes cuando se intentan rastrear las decisiones de los codificadores: Por ejemplo, ¿se dieron cuenta los codificadores de todos los informes de un evento, si todos los informes se refieren al mismo evento, si todos los informes se han procesado de forma coherente?, y así sucesivamente. Por lo tanto, la agregación por parte de los codificadores, sin la codificación por separado de las noticias, puede dificultar la determinación del origen de la baja concordancia entre codificadores en la inclusión y codificación de eventos.

### **Comparación de datos no mediáticos**

Para investigar los errores derivados de las preferencias de los medios de comunicación, varios investigadores han comparado los datos de eventos con fuentes de datos no mediáticas. Aunque estos datos y comparaciones son poco frecuentes, pueden dar indicaciones importantes sobre los errores de los medios de comunicación. Sin embargo, los propios datos externos pueden tener errores importantes (y desconocidos), lo que puede poner en peligro la validez de las afirmaciones, producidas de las comparaciones de los medios de comunicación.

Los registros policiales se utilizan con mayor frecuencia para investigar la cobertura mediática de eventos de protesta en contextos occidentales. Jenkins y Maher (2016:44) señalan que los estudios suelen afirmar que un solo periódico no cubre más (y a menudo menos) del 20 al 40% de los eventos identificados en los registros policiales. Aunque muchos estudios han utilizado los registros policiales para investigar el error de



cobertura (lo que confirma los efectos de selección analizados en la sección Cobertura de Noticias), observamos que también se han utilizado para estudiar el error de información. Nos referimos en particular al estudio de McCarthy et al. (1999:117-26) con respecto a los "hechos concretos" sobre las manifestaciones en Washington, DC (véase la sección de noticias).

No obstante, hay que tener cuidado para no confiar demasiado en la calidad de los registros policiales, ya que no se recogen necesariamente de forma sistemática y pueden carecer de detalles importantes de los eventos (Oliver y Myers 1999:48). En el caso de los eventos en contextos de países en vías de desarrollo -el foco geográfico de muchos conjuntos de datos sobre eventos conflictivos- los registros policiales pueden estar sujetos a errores más graves que en contextos occidentales, además de tener problemas de acceso (por ejemplo, Bocquier y Maupeu 2005:332).

Para los estudios centrados en los conflictos armados o en la violencia contra los civiles, los informes de las ONG son otra fuente externa y se utilizan comúnmente para la construcción de conjuntos de datos de eventos de conflicto (a menudo en adición a los datos de los medios de comunicación). Como muestran Davenport y Ball (2002) para la violencia estatal en Guatemala, los informes de las ONG documentan más violaciones del Estado y diferentes tendencias de la violencia estatal a lo largo del tiempo que los informes de los periódicos, aunque no puede establecerse si esto se debe a un error de medición o de representación. Curiosamente, los datos de las entrevistas muestran otro panorama. Además, aunque el propósito de muchas ONG sobre el terreno es proporcionar información independiente y fiable, la información puede depender de la atención de los donantes a los eventos "candentes" o ser creada deliberadamente para atraer la atención de los medios de comunicación internacionales. A su vez, los informes de las ONG se basan a menudo en los medios de comunicación. Por lo tanto, los informes de las ONG podrían reforzar el sesgo de los medios de comunicación hacia determinados países o conflictos en los conjuntos de datos sobre eventos. Aunque un set de datos militares puede revelar importantes conocimientos sobre la cobertura de los conflictos armados (Weidmann 2015, 2016), también puede servir a determinados objetivos organizativos y no ser necesariamente un reflejo fiel de la realidad.

### Ajustes de los datos

Demarest, Leila y Langer, Arnim "Cómo se incorporan (o no) los eventos en los sets de datos: Los obstáculos y las orientaciones en la utilización de los periódicos en el estudio de los conflictos" *Revista de Estudios Marítimos y Sociales*, N°24, ene-jun 2024, pp. 197-244.





## **Ponderación de los datos.**

El último paso del proceso de investigación de eventos es la fase de análisis, durante la cual los investigadores pueden aplicar correcciones al conjunto de datos de eventos para compensar los errores de muestreo o de medición. Un primer tipo de corrección consiste en ponderar los datos para corregir la subrepresentación de determinados eventos. Aunque la intención es reducir el error, este tipo de corrección también puede crearlo. De hecho, a menudo no hay datos externos que coincidan con el set de datos de eventos basados en los medios de comunicación. Entonces se hacen correcciones basadas en diferentes estudios y se supone que estas afirmaciones se mantienen en el espacio y el tiempo. Hug y Wisler (1998), por ejemplo, proponen correcciones estadísticas para el sesgo de selección (por ejemplo, la ponderación) basadas en un estudio comparativo de los registros policiales y los periódicos locales de cuatro ciudades suizas. Sostienen que las correcciones para las preferencias de cobertura de eventos violentos y eventos con más participantes también podrían ser útiles en otros contextos. Ortiz et al. (2005:408-11) sostienen por el contrario que éste puede ser un procedimiento que aumente el sesgo si la relevancia de los factores de selección, así como su magnitud, no se traslada a otros contextos.

Otras correcciones propuestas recientemente no se basan en comparaciones con fuentes externas, sino con los datos de otros medios. Hendrix y Salehyan (2015) utilizan un método de marcado y recaptura para estimar el verdadero número de eventos basándose en la información de múltiples fuentes de medios de comunicación. La SCAD recurre a informes de Associated Press (AP) y Agence France-Presse (AFP). El esquema de codificación, a partir de 2012, registra si un evento fue reportado en AFP, AP, o en ambas fuentes. Al estimar la correspondencia entre las fuentes, es posible hacer correcciones a los datos para los eventos no cubiertos en ambos conjuntos de datos. Se utiliza un enfoque similar propuesto por Cook et al. (2017). Es importante destacar que el método requiere que los conjuntos de datos informen sistemáticamente de todas las fuentes que han informado sobre un evento, lo que no es una práctica habitual. De hecho, aunque los sets de datos suelen citar una fuente concreta, esto no implica que el evento no haya quedado incluido en otras fuentes. El SCAD es una notable excepción. Sin embargo, se basa en el mismo tipo de fuentes de medios de comunicación, los

Demarest, Leila y Langer, Arnim "Cómo se incorporan (o no) los eventos en los sets de datos: Los obstáculos y las orientaciones en la utilización de los periódicos en el estudio de los conflictos" *Revista de Estudios Marítimos y Sociales*, N°24, ene-jun 2024, pp. 197-244.





servicios de noticias internacionales, mientras que los periódicos locales podrían captar un número considerable de eventos adicionales (véase la sección de Cobertura de Datos). Para corregir la atención diferencial hacia determinados países por parte de los medios de comunicación internacionales (por ejemplo, Herkenrath y Knoll 2011), también se ha propuesto incluir una variable para el número total de noticias no relacionadas con el conflicto dedicadas a un país concreto en un año determinado como variable de control en los análisis sustantivos (Hendrix y Salehyan 2017:1664-65).

### **Imputación de datos perdidos.**

Un segundo tipo de corrección que puede hacerse a los datos de eventos conflictivos es la imputación de los datos que faltan para compensar errores de medición. Estos ajustes pueden, a su vez, dar lugar a estadísticas erróneas. A menudo se realizan correcciones de datos ausentes para las fechas, las geolocalizaciones y las estimaciones de mortalidad. Por ejemplo, el UCDP GED da una fecha y una hora a cada evento, pero para algunos eventos, surge la incertidumbre sobre la precisión de estas variables (Croicu y Sundberg 2016:5-6). A veces solo se conoce la semana, el mes o el año de un evento. En estos casos, UCDP GED asigna la fecha más temprana posible al evento. También es habitual dar las coordenadas geográficas del centro de la unidad administrativa o del país cuando se desconoce la ubicación exacta. La imputación de datos temporales y de localización suele ir acompañada de variables que indican un nivel de incertidumbre en la codificación. ACLED adopta enfoques similares (Raleigh y Dowd 2017:14-16). Es importante destacar que no está claro hasta qué punto los indicadores de precisión se utilizan realmente en las aplicaciones empíricas de los datos de eventos de conflicto, por ejemplo, excluyendo los eventos inciertos como una comprobación de robustez.

Por último, también existen imputaciones para las estimaciones de la mortalidad. Un ejemplo es la división del recuento de víctimas cuando un evento ocurrió en varios lugares o en el transcurso de varias fechas (por ejemplo, ACLED pero no SCAD). Otro ejemplo se refiere a las palabras utilizadas para describir el número de víctimas, lo que es relativamente común (por ejemplo, varios, algunos, docenas). Chojnacki et al. (2012:391-92) optan por anotar la palabra en el conjunto de datos pero sin cuantificarla. SCAD (Codebook 3.1., actualizado el 20 de noviembre de 2017:5) hace uso de una

Demarest, Leila y Langer, Arnim "Cómo se incorporan (o no) los eventos en los sets de datos: Los obstáculos y las orientaciones en la utilización de los periódicos en el estudio de los conflictos" *Revista de Estudios Marítimos y Sociales*, N°24, ene-jun 2024, pp. 197-244.





distinción para desaparecidos pero más ("probablemente grande") o menos ("probablemente pequeño") que 10. ACLED (Raleigh y Dowd 2017:20) opta por cuantificar la descripción: Varios, muchos, plural o desconocido se establece en 10, docenas se establece en 12, cientos se establece en 100. Esta cuantificación podría poner en peligro la calidad de los datos.

### **Datos de eventos: Un camino a seguir**

El entorno (framework) TEE esbozado en la sección anterior ha permitido un debate exhaustivo sobre las causas de error que pueden afectar a la calidad de los datos de los eventos conflictivos, que abarcan todas las subdisciplinas de especialización. También hemos discutido las posibles estrategias para mitigar estos errores propuestas en la literatura, así como sus límites. El cuadro 1 ofrece una visión general de estas fuentes de error, las estimaciones disponibles de su tamaño y estrategias para reducirlos. Las estimaciones del grado de error se basan en los estudios aquí revisados y, por tanto, en diferentes contextos geográficos, períodos de tiempo, procedimientos de codificación (automatizados), etc. Además, las estimaciones muestran hasta qué punto puede desviarse la información, pero no necesariamente cómo repercute esto en los resultados sustantivos de la investigación. Aunque esto debe tenerse en cuenta, proporcionan a los investigadores indicaciones sobre cómo evaluar la calidad de los datos. Por último, las estimaciones disponibles (y no disponibles) también indican dónde se necesita más investigación empírica. En esta sección, nos centramos sobre todo en las preguntas metodológicas que hasta ahora no se han abordado suficientemente en la bibliografía. La última columna del cuadro 1 contiene una extensa lista de preguntas que requieren más investigación y que, en conjunto, constituyen una agenda de investigación relativa a la metodología de creación y uso de los conjuntos de datos de eventos conflictivos.

La mayor parte de la atención en la literatura se ha dirigido al error de cobertura y por razones importantes. De hecho, las estimaciones disponibles sobre la selección de eventos revelan que los efectos de distorsión del error de cobertura en los resultados de la investigación pueden ser sustanciales. Aunque la mayor parte de las pruebas de este sesgo se han establecido en el contexto de los movimientos de protesta en occidente (por ejemplo, Jenkins y Maher 2016), hay indicios de que los conflictos violentos también son reportados de manera insatisfactoria (Davenport y Ball 2002, Weidmann

Demarest, Leila y Langer, Arnim "Cómo se incorporan (o no) los eventos en los sets de datos: Los obstáculos y las orientaciones en la utilización de los periódicos en el estudio de los conflictos" *Revista de Estudios Marítimos y Sociales*, N°24, ene-jun 2024, pp. 197-244.





2016). Además de la forma del conflicto otro reto importante es el uso generalizado de noticias internacionales para investigar los conflictos (violentos) en el mundo en desarrollo. La evidencia sugiere que esto puede ser problemático (Bueno de Mesquita et al. 2015, Demarest y Langer 2018, Herkenrath y Knoll 2011). Este problema podría mitigarse con la mayor disponibilidad de fuentes locales en línea y, potencialmente, con nuevas evoluciones en la codificación automatizada. Curiosamente, las estimaciones disponibles sobre el error de reporte parecen indicar que los hechos de las protestas (McCarthy et al. 1999) y los conflictos violentos (Weidmann 2015) pueden ser reportados relativamente libres de errores. No obstante, el error de información también puede depender del contexto. Esto es especialmente importante a tener en cuenta cuando se trata de medios de comunicación locales sometidos al control del gobierno. Sin embargo, no parece haber estimaciones disponibles para este tipo de contextos.

Los errores derivados de la lógica de los medios de comunicación requieren una cuidadosa consideración y una mayor investigación. Otras características del protocolo de recogida de datos también requieren atención. Así lo revelan también las estimaciones de los errores relacionados con el proceso de codificación. A este respecto, es importante señalar que, si bien se pueden utilizar diferentes indicadores para evaluar determinadas opciones metodológicas (acuerdo de selección, fiabilidad entre codificadores, recuperación y precisión), por el momento no existe un consenso real en la bibliografía sobre el uso de dichos indicadores y, en consecuencia, su presentación. Por ejemplo, muchos de los nuevos conjuntos de datos de los estudios sobre la paz y los conflictos, rara vez proporcionan información sobre la selección y fiabilidad de los codificadores o el grado general de datos imputados en el conjunto de datos. En cambio, los desarrolladores de codificación automatizada parecen estar más de acuerdo en la necesidad de informar sobre los índices de recuperación y precisión.

La medición de estos errores es importante para establecer hacia dónde deben dirigirse los esfuerzos de recopilación de datos para lograr los mayores beneficios en términos de calidad de los datos. Por ejemplo, Schrodtt y Ulfelder (2016:29) sostienen que incluir indicadores de incertidumbre sobre eventos, actores, etc. (véase el cuadro 1, "desarrollo del libro de códigos"), no añadió mucho valor a los datos de atrocidades de PITF, y dejaron fuera estos indicadores en versiones posteriores. Las mismas preguntas pueden

Demarest, Leila y Langer, Arnim "Cómo se incorporan (o no) los eventos en los sets de datos: Los obstáculos y las orientaciones en la utilización de los periódicos en el estudio de los conflictos" *Revista de Estudios Marítimos y Sociales*, N°24, ene-jun 2024, pp. 197-244.







plantearse con respecto a la codificación de los informes por separado para tener en cuenta las diferencias entre ellos (Day et al. 2015, Weidmann y Rød 2015). Es necesario seguir investigando para determinar las ventajas de estos procedimientos y poder evaluar su uso en nuevos conjuntos de datos.

Los usuarios también deben prestar suficiente atención a las fuentes de error de los eventos. Esto se aplica a la selección de determinados sets de datos para abordar preguntas de investigación sustantivas, pero también en la presentación de informes y en las comprobaciones de robustez. Las causas de los errores de evento son necesarias para comprender los límites de determinados estudios, por ejemplo, tanto en el ámbito académico como en el político. Además, cuando los desarrolladores de datos proporcionan indicadores de la calidad de los mismos, sostenemos que los investigadores que se centran en preguntas sustantivas no sólo deberían informar de estos indicadores en sus estudios, sino que también deberían reflexionar sobre las implicaciones de estos indicadores para la validez de sus afirmaciones y conclusiones. Esto incluye, en particular, los indicadores de imputación de datos faltantes, que no se utilizan habitualmente en los estudios cuantitativos sobre conflictos, a pesar de que este ámbito se centra cada vez más en detalles precisos sobre los eventos, tanto en el tiempo como en el lugar (Gleditsch et al. 2014). Aunque la ponderación de los datos de eventos todavía no se utiliza habitualmente, de nuevo el efecto de la ponderación debe compararse cuidadosamente con los resultados basados en datos no ponderados, y debe motivarse explícitamente la preferencia por unos resultados sobre otros.



Tabla 1. Errores totales en eventos: errores, orientaciones e investigación futura

Etapa de investigación	Riesgos de errores	Estimaciones disponibles <sup>5</sup>	Soluciones	Futuras líneas de investigación
Muestreo de fuentes de noticias	<p>Representación y medición:</p> <p>La cobertura informativa y los reportajes dependen de:</p> <ul style="list-style-type: none"> <li>-Características del evento: selección más probable de eventos violentos, eventos con un mayor número de participantes/información de las protestas más fácil de constatar que de los eventos de conflictos armados (por ejemplo, ubicación).</li> <li>-Características de la fuente de noticias: ideología, escala geográfica del público destinatario (local, nacional, internacional),</li> </ul>	<p><u>Representación:</u></p> <ul style="list-style-type: none"> <li>-Los informes de las noticias nacionales frente a los registros policiales<sup>6</sup>: 20-60 por ciento de los eventos de protesta.</li> <li>-Las noticias nacionales registran alrededor de 10 veces menos violencia letal que las fuentes documentales de las organizaciones no gubernamentales (ONG) y unas 2 veces menos que las entrevistas<sup>7</sup>.</li> <li>-Noticias internacionales frente a datos militares<sup>8</sup>: 28,5% de eventos letales.</li> <li>-Noticias internacionales frente a nacionales<sup>9</sup>: 1,5 -5 por ciento de protestas/disturbios; alrededor del 25 por ciento de eventos mortales; alrededor de 0,7 veces el número de eventos terroristas.</li> <li>-Noticias nacionales frente a provinciales<sup>10</sup>: alrededor de 0,3 veces el número de muertes.</li> </ul> <p><u>Mediciones:</u></p> <ul style="list-style-type: none"> <li>-Correspondencia del informe de protesta con los registros policiales<sup>11</sup>:</li> </ul>	<ul style="list-style-type: none"> <li>-Extraer datos de múltiples fuentes de noticias, incluidas fuentes de los medios de comunicación (con diferentes orientaciones políticas), así como fuentes externas (por ejemplo, informes de ONG).</li> <li>-Utilice fuentes de noticias nacionales o locales para los estudios de un solo país. Las fuentes internacionales pueden ser necesarias para estudios transnacionales a gran escala, pero conllevan riesgos de representación.</li> <li>-Adaptar la escala de los eventos conflictivos estudiados a la escala de la fuente de noticias. Las fuentes locales pueden ser más adecuadas para estudiar eventos de bajo nivel (por ejemplo, protestas frente a atentados terroristas) y realizar análisis subnacionales.</li> <li>-Aplicar ponderaciones correctoras para compensar los efectos de la cobertura basándose en comparaciones entre fuentes de medios de comunicación o con datos externos (por ejemplo, registros policiales).</li> <li>-Control de la cantidad de informes no</li> </ul>	<ul style="list-style-type: none"> <li>-¿Qué eventos conflictivos tienen más probabilidades de aparecer en las noticias? ¿Cómo difieren los sesgos según el tipo de eventos (por ejemplo, protestas frente a conflictos violentos)?</li> <li>-¿Cómo difieren la cobertura y la información entre las fuentes de noticias locales, nacionales e internacionales?</li> <li>-¿Cómo podemos utilizar procedimientos automatizados de codificación en las fuentes de noticias locales teniendo en cuenta el acceso a los</li> </ul>

<sup>5</sup> Los porcentajes se utilizan cuando los eventos/fatalidades coinciden; en caso contrario, se calcula la diferencia en el número de eventos/fatalidades registrados (X veces menos o más).

<sup>6</sup> Jenkins y Maher (2016); Myers y Canigla (2004).

<sup>7</sup> Davenport y Ball (2022).

<sup>8</sup> Weidmann (2016)

<sup>9</sup> Demarest and Langer (2018), Herkenrath y Knoll (2011), Bueno de Mesquita et al. (2015)

<sup>10</sup> Barron and Sharpe (2008)

Demarest, Leila y Langer, Arnim "Cómo se incorporan (o no) los eventos en los sets de datos: Los obstáculos y las orientaciones en la utilización de los periódicos en el estudio de los conflictos" *Revista de Estudios Marítimos y Sociales*, N°24, ene-jun 2024, pp. 197-244.



Continuación tabla 1

Etapa de investigación	Riesgos de errores	Estimaciones disponibles	Soluciones	Futuras líneas de investigación
	preferencia de la fuente (por ejemplo, fuentes gubernamentales). -Características del contexto: por ejemplo, infraestructura deficiente para acceder a la información (o verificarla), control gubernamental de la información.	>,98(fecha),>,65(propósito),>61(tamaño). -Correspondencia de informes sobre conflictos armados con datos militares <sup>12</sup> : -80 por ciento a menos de 50 km de la ubicación real; 50 por ciento de bajas por correspondencia, pequeñas diferencias	relacionados con el conflicto para tener en cuenta las diferencias en la atención de los medios de comunicación hacia países o regiones. -Codifique explícitamente las diferencias entre las fuentes y téngalas en cuenta en los análisis sustantivos para tener en cuenta los errores de los informes.	datos digitales, las diferencias lingüísticas, etc.? -¿Cómo afecta la orientación política de la fuente de noticias a la cobertura e información de los eventos? -¿Cómo afecta la libertad de prensa a la cobertura de los eventos?
Muestreo de fuentes de noticias	<i>Representación:</i> -El muestreo por ejemplar (por ejemplo, los números del lunes) y por página (por ejemplo, la primera página) aumenta la falta de fiabilidad y puede reforzar el sesgo de cobertura. <i>Medición:</i> -Los informes pueden basarse en fuentes sesgadas (por ejemplo, gobierno, rebeldes, partidos políticos). -En los informes falta información sobre la fecha, el lugar, los actores y las víctimas. -Los informes ponen	-11-18 por ciento de eventos sin tiempo preciso, 23-49 por ciento sin lugar preciso <sup>13</sup> .  -Los indicadores de incertidumbre marcan poca diferencia <sup>14</sup> .	-Codifique todos los números y páginas de un subconjunto (aleatorio) de los datos y compárelos con la muestra reducida. Posible aplicación de ponderaciones correctoras. -Incluir indicadores de falta de fiabilidad y sesgo en el libro de códigos para la ocurrencia de eventos, identidades de los actores, estimaciones de víctimas mortales, etc. -Codificar por separado los informes sobre un mismo evento e incluir un rango de ambigüedad en el set de datos final o permitir	-¿Cuáles son los efectos de los distintos mecanismos de muestreo de informes en las estadísticas de eventos (por ejemplo, muestreo de la primera página, muestreo del número completo, cadena de búsqueda en el repositorio en línea, procedimiento automatizado)? - ¿Cómo afecta la inclusión de indicadores de falta de fiabilidad y sesgo a las afirmaciones sustantivas de los análisis de conflictos?

<sup>11</sup> McCarthy et al. (1999), estimaciones de los medios impresos.

<sup>12</sup> Weidmann (2015)

<sup>13</sup> Datos de ACLED y UCDP, respectivamente.

Demarest, Leila y Langer, Arnim "Cómo se incorporan (o no) los eventos en los sets de datos: Los obstáculos y las orientaciones en la utilización de los periódicos en el estudio de los conflictos" *Revista de Estudios Marítimos y Sociales*, N°24, ene-jun 2024, pp. 197-244.





Continuación tabla 1

Etapa de investigación	Riesgos de errores	Estimaciones disponibles	Soluciones	Futuras líneas de investigación
	explícitamente en duda la ocurrencia de los eventos, la identidad de los actores y las estimaciones de víctimas mortales.		a los usuarios probar diferentes especificaciones (por ejemplo, mínimo y máximo de víctimas mortales).	
Desarrollo del libro de códigos	<p><i>Representación/Medición:</i></p> <ul style="list-style-type: none"> <li>- Las instrucciones de selección/codificación poco claras y/o ambiguas pueden crear incoherencias en el proceso de codificación y diferencias tanto asistemáticas como sistemáticas entre los codificadores.</li> <li>- Los diccionarios para la codificación automatizada pueden dejar fuera demasiados eventos relevantes e introducir demasiados irrelevantes y codificar variables de forma incorrecta.</li> </ul>	<ul style="list-style-type: none"> <li>- Muestreo de eventos relevantes<sup>15</sup>: 80-97% de recuperación y 23-58% de precisión.</li> <li>- Categorización correcta de los eventos<sup>16</sup>: 7-96 por ciento.</li> </ul>	<ul style="list-style-type: none"> <li>- Calcular las medidas de selección/fiabilidad entre codificadores humanos en una fase piloto y a lo largo del proceso de codificación para adaptar las instrucciones cuando sea necesario.</li> <li>- Ajustar el diccionario para los procedimientos automatizados tras la comprobación de los datos y volver a ejecutar el programa de codificación.</li> </ul>	<ul style="list-style-type: none"> <li>- ¿Cómo afectan las instrucciones del libro de códigos a la fiabilidad de la selección/codificación entre codificadores?</li> <li>- ¿Qué tipo de instrucciones funcionan mejor teniendo en cuenta la realidad de los datos (por ejemplo, informes imprecisos)?</li> <li>- ¿Qué nivel de complejidad podemos alcanzar con la codificación humana y automatizada?</li> <li>- ¿Cuáles son las posibilidades de mejorar y ampliar la codificación automatizada a otros idiomas, contextos, etc.?</li> </ul>

<sup>14</sup> Schrodtt y Ulfelder (2016).

<sup>15</sup> Bond et al. (1997), Croicu y Weidmann (2016), King y Lowe (2003).

<sup>16</sup> Bond et al. (1997), Boschee et al. (2013), King y Lowe (2003), Stepinski et al. (2006).

Demarest, Leila y Langer, Arnim "Cómo se incorporan (o no) los eventos en los sets de datos: Los obstáculos y las orientaciones en la utilización de los periódicos en el estudio de los conflictos" *Revista de Estudios Marítimos y Sociales*, Nº24, ene-jun 2024, pp. 197-244.





Continuación tabla 1

Etapa de investigación	Riesgos de errores	Estimaciones disponibles	Soluciones	Futuras líneas de investigación
Proceso de codificación	<p><i>Representación:</i></p> <ul style="list-style-type: none"> <li>- Los codificadores omiten eventos relevantes o incluyen eventos irrelevantes.</li> <li>- Provoca falta de fiabilidad, pero también un sesgo potencial cuando los codificadores pasan por alto eventos específicos con más frecuencia (por ejemplo, protestas de perfil bajo).</li> </ul> <p><i>Medición:</i></p> <ul style="list-style-type: none"> <li>- Los codificadores clasifican erróneamente eventos o actores, cometen errores en los datos de tiempo y lugar, etc.</li> <li>- Puede causar falta de fiabilidad o sesgo cuando los codificadores malinterpretan sistemáticamente las instrucciones</li> </ul>	<ul style="list-style-type: none"> <li>- Solapamiento del 60-70 por ciento entre codificadores independientes en la selección de eventos de protesta.<sup>17</sup></li> <li>- Las ks de Cohen muestran una concordancia sustancial para la categoría de evento, actores mortales, causa<sup>18</sup>.</li> <li>- Las alpha de Krippendorff a para categoría de evento, actores, víctimas mortales &lt;0,77.<sup>19</sup></li> </ul>	<ul style="list-style-type: none"> <li>- Seguimiento del proceso de codificación.</li> <li>- Calcule medidas de selección/fiabilidad entre codificadores en una fase piloto y a lo largo de todo el proceso de codificación.</li> <li>- Retener a los codificadores que detecten los eventos más relevantes/tengan puntuaciones de fiabilidad más altas tras una fase inicial de prueba.</li> <li>- Haga que un equipo de investigación formado por separado preseleccione los eventos relevantes que serán codificados por un equipo diferente.</li> </ul>	<ul style="list-style-type: none"> <li>- ¿Cuáles son las medidas de selección/fiabilidad de los intercodificadores para los sets de datos conocidos actualmente?</li> <li>- ¿Es mayor la fiabilidad de la selección cuando se utilizan cadenas de búsqueda o procedimientos semiautomáticos?</li> <li>- ¿Cómo interactúa la selección con las características del evento (por ejemplo, la violencia) y la fuente de noticias (por ejemplo, el uso habitual de visualizaciones)?</li> <li>- ¿Hasta qué punto es fiable la codificación de hechos concretos (por ejemplo, fecha, ubicación, actor) frente a hechos no concretos (por ejemplo, causa atribuida del evento)?</li> <li>- ¿Cuáles son las características de los "buenos codificadores" (por ejemplo, nivel de estudios, duración del contrato)?</li> </ul>

<sup>17</sup> Hutter (2014), Kriesi et al. (1998).

<sup>18</sup> Salehyan et al. (2012)

<sup>19</sup> Demarest y Langer (2018)

Demarest, Leila y Langer, Arnim "Cómo se incorporan (o no) los eventos en los sets de datos: Los obstáculos y las orientaciones en la utilización de los periódicos en el estudio de los conflictos" *Revista de Estudios Marítimos y Sociales*, N°24, ene-jun 2024, pp. 197-244.





Continuación tabla 1

Etapa de investigación	Riesgos de errores	Estimaciones disponibles	Soluciones	Futuras líneas de investigación
Comparación de datos no procedentes de los medios de comunicación	<p><i>Representación y Medición:</i></p> <ul style="list-style-type: none"> <li>- Los datos externos (por ejemplo, registros policiales, datos militares, informes de ONG) también pueden adolecer de sesgos y falta de fiabilidad.</li> <li>- Cuando se utilizan para corregir datos basados en los medios de comunicación, los errores en los datos externos pueden crear una incertidumbre adicional o reforzar los sesgos.</li> </ul>	<ul style="list-style-type: none"> <li>- Los informes de las ONG contienen unas 4 veces más muertes que los datos de las entrevistas<sup>20</sup>.</li> </ul>	<ul style="list-style-type: none"> <li>- Comparar diferentes fuentes no mediáticas (por ejemplo, informes de ONG, entrevistas, encuestas) para investigar la cobertura específica de cada fuente y los efectos de la información.</li> <li>- Recurrir a los medios de comunicación y a fuentes externas, codificar los informes por separado e incluir indicadores de posible incertidumbre y sesgo en el set de datos final.</li> </ul>	<ul style="list-style-type: none"> <li>- ¿Qué decisiones toman los redactores de los informes de las ONG y cómo puede afectar esto a los datos de los eventos (por ejemplo, es probable que exageren las atrocidades para abogar por más apoyo)?</li> <li>- ¿Hasta qué punto podemos confiar en los registros policiales/datos militares en contextos en desarrollo?</li> </ul>
Ajuste de datos	<p><i>Representación:</i></p> <ul style="list-style-type: none"> <li>- Las ponderaciones correctoras para compensar los efectos de selección pueden inducir sesgos cuando el sesgo de selección no es constante en el tiempo y el espacio.</li> </ul> <p><i>Medición:</i></p> <ul style="list-style-type: none"> <li>- La imputación de los datos que faltan sobre fechas, lugares, víctimas mortales, etc., puede dar lugar a una falsa sensación de fiabilidad.</li> <li>- La imputación temporal y de datos puede afectar, en particular, a los análisis que requieren datos temporales y de localización precisos.</li> </ul>	<ul style="list-style-type: none"> <li>- Sesgo en la significación y la dirección de los coeficientes de regresión.<sup>21</sup></li> </ul>	<ul style="list-style-type: none"> <li>- Crear ponderaciones en el set de datos, garantizar la transparencia en su creación y permitir a los usuarios incorporar o no la ponderación.</li> <li>- Incluir indicadores para la imputación de variables.</li> </ul>	<ul style="list-style-type: none"> <li>- ¿Cómo afectan las medidas correctoras a las afirmaciones sustantivas?</li> <li>- ¿Cómo varía el sesgo de selección en función del tiempo, el espacio y la fuente de noticias?</li> <li>- ¿Cómo afecta la imputación de datos a las conclusiones sustantivas?</li> </ul>

<sup>20</sup> Davenport y Ball (2002)

<sup>21</sup> Ortiz et al. (2005)

Demarest, Leila y Langer, Arnim "Cómo se incorporan (o no) los eventos en los sets de datos: Los obstáculos y las orientaciones en la utilización de los periódicos en el estudio de los conflictos" *Revista de Estudios Marítimos y Sociales*, N°24, ene-jun 2024, pp. 197-244.



## Conclusión

La calidad de los datos sobre eventos conflictivos puede verse afectada por una amplia gama de errores. El debate de este artículo se ha guiado por el estado del arte actual sobre los estudios de eventos conflictivos y también se basó en los estudios sobre movimientos sociales y comunicación. La principal ventaja del entorno TEE es que ofrece una perspectiva holística de las causas de error que afectan a los datos de los conflictos y, en consecuencia, claridad analítica en un campo de estudio discutiblemente amplio. De hecho, aunque se han debatido muchas causas de error en la literatura, estos debates no siempre han permitido una mayor sistematización. Al hacer justamente esto, el entorno TEE ofrece una herramienta de referencia para los desarrolladores y usuarios de datos, tanto nuevos como veteranos, así como una guía para futuras investigaciones. Además, aunque TEE se ha centrado en las prácticas de recogida de datos de eventos humanos y automatizados, puede ampliarse a nuevas áreas. La aparición de la "información ciudadana" a través de las redes sociales, por ejemplo, se está convirtiendo en una nueva e importante fuente de recopilación de datos de eventos, pero se plantean preocupaciones similares con respecto a los efectos de la cobertura y la información, así como a los procedimientos de recopilación de datos.

Por último, cabe señalar que, aunque los errores pueden y deben minimizarse, difícilmente puedan descartarse por completo. Por lo tanto, los datos de eventos nunca serán un reflejo fiel de la realidad. Sin embargo, esto no es diferente de otras fuentes de datos empíricos en las ciencias sociales, incluyendo las encuestas de opinión pública. Volviendo a nuestra analogía inicial, vale la pena considerar que las fuentes de error están ampliamente reconocidas en la investigación con encuestas, pero también que el verdadero ejercicio consiste en minimizar los errores teniendo en cuenta los recursos limitados. Al igual que ocurre con "los errores y los costes de las encuestas" (Groves 1989), el equilibrio entre los errores y los costes de los eventos limita a los desarrolladores de conjuntos de datos sobre eventos. Sin embargo, para mejorar las directrices y las normas de recopilación de datos, la metodología de los datos de eventos conflictivos debe considerarse un programa de investigación por derecho propio. El marco TEE y las preguntas de investigación expuestas en el cuadro 1 ofrecen importantes orientaciones para ello.



## Notas

1. Las versiones anteriores del conjunto de datos de eventos georreferenciados del UCDP sólo incluían eventos de conflicto en África, actualmente está disponible un nuevo set de datos globales (versión 5.0). Obsérvese que, especialmente en el caso de los conjuntos de datos de eventos de conflictos violentos, predominan los países en desarrollo, aunque el set de datos tenga un enfoque global.
2. Se ha escrito una amplia gama de estudios sobre el tema del sesgo de selección de los medios de comunicación. Aquí ofrecemos una visión general de las ideas clave y dirigimos a los lectores, para obtener más información a las referencias citadas en esta sección, y a la sección de ponderación de datos para las correcciones sobre el sesgo de selección.
3. Todos los conjuntos de datos sobre eventos conflictivos mencionados en la introducción utilizan predominantemente estos inventarios, excepto el Conjunto de Datos de Localización y Eventos de Conflictos Armados (en inglés, Armed Conflict Location and Event Dataset), que también se basa en los periódicos locales.
4. Sin embargo, el uso de fuentes de noticias locales ha sido más frecuente para investigar los disturbios entre hindúes y musulmanes en la India (por ejemplo, Wilkinson 2004).
5. En la literatura, las noticias duras también se conceptualizan como las que tienen un alto grado de interés informativo, como las noticias sobre política, economía y asuntos sociales de interés periodístico, mientras que las noticias blandas tienen un valor informativo menos sustantivo, por ejemplo chismes, historias de interés humano, etc. (por ejemplo, Reinemann et al. 2011). Esta distinción está relacionada, pero difiere de la utilizada en este artículo.
6. Aunque el diccionario y el programa de codificación son en principio entidades separadas en procedimientos de codificación automatizada (Schrodt y Van Brackle 2013:24), no separamos explícitamente a los dos en nuestra discusión aquí. Tampoco profundizamos en los errores o características de la programación (por ejemplo, la velocidad).





7. La recuperación (recall) es igual al número de verdaderos positivos dividido por la suma del número de verdaderos positivos y el número de falsos negativos. La precisión es igual al número de verdaderos positivos dividido por la suma del número de verdaderos positivos y el número de falsos positivos. La estadística F1 captura la media armónica de recuperación y precisión (Heap et al. 2017).
8. El procedimiento de análisis sintáctico disperso desglosa las frases del texto relevante en función de actores, objetivos y verbos transitorios.
9. Sin embargo, es importante mencionar que el uso de TABARI en su trabajo ha sido criticado por Schrodt (Schrodt y Van Brackle 2013:27).
10. Véase el Libro de códigos 3.1, actualizado el 20 de noviembre de 2014, pp. 3-4.
11. Además de la fiabilidad entre codificadores, se puede prestar atención a la fiabilidad de intracodificadores o la estabilidad, es decir, si un codificador codifica informes anteriores de la misma manera que en un momento posterior (Krippendorff 2013:270-71).

## Bibliografía

### **AZAR, EDWARD E.**

1980 The Conflict and Peace Data Bank (COPDAB) Project. *Journal of Conflict Resolution* 24:143-52.

### **BARRON, PATRICK AND JOANNE SHARPE.**

2008 Local Conflict in Post-Suharto Indonesia: Understanding Variations in Violence Levels and Forms through Local Newspapers. *Journal of East Asian Studies* 8:395-423.

### **BERNAUER, THOMAS AND NILS P. GLEDITSCH.**

2012 New Event Data in Conflict Research. *International Interactions* 38:375-81.

### **BOCQUIER, PHILIPPE AND HERVÉ MAUPEU.**

2005 Analysing Low Intensity Conflict in Africa Using Press Reports. *European Journal of Population* 21:321-45.

### **BOND, DOUG, JOE BOND, CHURL OH, J. CRAIG JENKINS, AND CHARLES LEWIS TAYLOR.**

2003 Integrated Data for Events Analysis (IDEA): An Event Typology for Automated Events Data Development. *Journal of Peace Research* 40:733-45.

### **BOND, DOUG, J. CRAIG JENKINS, CHARLES L. TAYLOR, AND KURT SCHOCK.** 1997 Mapping Mass Political Conflict and Civil Society: Issues and Prospects for the



Automated Development of Event Data. *Journal of Conflict Resolution* 41:553-79.

**BOSCHEE, ELIZABETH, PREMKUMAR NATARAJAN, AND RALPH WEISCHEDEL.** 2013 Automatic Extraction of Events from Open Source Text for Predictive Forecasting in *Handbook of Computational Approaches to Counterterrorism*, edited by V. S. Subrahmanian. New York: Springer Science & Business Media: 51-67

**BUENO DE MESQUITA, ETHAN, C. CHRISTINE FAIR, JENNA JORDAN, RASUL B. RAIS, AND JACOB N. SHAPIRO.**

2015 Measuring Political Violence in Pakistan: Insights from the BFRS Dataset. *Conflict Management and Peace Science* 32:536-58.

**CHAN, JOSEPH M. AND CHI-CHUAN LEE.**

1984 Journalistic 'Paradigms' of Civil Protests: A Case Study in Hong Kong, in *The News Media in National and International Conflict*, edited by Andrew Arno and Wimal Dissanayake. Boulder, CO: Westview: 249-76

**CHOJNACKI, SVEN, CHRISTIAN ICKLER, MICHAEL SPIES, AND JOHN WIESEL.** 2012 Event Data on Armed Conflict and Security: New Perspectives, Old Challenges, and Some Solutions. *International Interactions* 38:382-401.

**COLLIER, PAUL AND ANKE HOFFFLER.**

2002 Greed and Grievance in Civil War. *Centre for the Study of African Economies Working Paper Series 2002-01*, Oxford, England.

**COOK, SCOTT J., BETSABE BLAS, RAYMOND J. CARROLL, AND SAMIRAN SINHA.**

2017 Two Wrongs Make a Right: Addressing Underreporting in Binary Data from Multiple Sources. *Political Analysis* 25:223-40.

**CROICU, MIHAI AND RALPH SUNDBERG.**

2016 UCDP GED Codebook Version 5.0. Uppsala, Sweden: Department of Peace and Conflict Research, Uppsala University.

**CROICU, MIHAI AND NILS B. WEIDMANN.**

2015 Improving the Selection of News Reports for Event Coding Using Ensemble Classification. *Research & Politics* 2:1-8.

**DARDIS, FRANK E.**

2006 Marginalization Devices in U.S. Press Coverage of Iraq War Protest: A Content Analysis. *Mass Communication and Society* 9:117-35.

**DANZGER, M. HERBERT.**

1975 Validating Conflict Data. *American Sociological Review* 40:570-84.

**DAVENPORT, CHRISTIAN.**

2010 *Media Bias, Perspective, and State Repression: The Black Panther Party*. Cambridge, MA: Cambridge University Press.

**DAVENPORT, CHRISTIAN AND PATRICK BALL.**

Demarest, Leila y Langer, Arnim "Cómo se incorporan (o no) los eventos en los sets de datos: Los obstáculos y las orientaciones en la utilización de los periódicos en el estudio de los conflictos" *Revista de Estudios Marítimos y Sociales*, N°24, ene-jun 2024, pp. 197-244.



2002 Views to a Kill: Exploring the Implications of Source Selection in the Case of Guatemalan State Terror, 1977-1995. *Journal of Conflict Resolution* 46:427-50.

**DAY, JOEL, JONATHAN PINCKNEY, AND ERICA CHENOWETH.**

2015 Collecting Data on Nonviolent Action: Lessons Learned and Ways Forward. *Journal of Peace Research* 52:129-33.

**DEMAREST, LEILA AND ARNIM LANGER.**

2018 The Study of Violence and Social Unrest in Africa: A Comparative Analysis of Three Conflict Event Datasets. *African Affairs* 117:310-25.

**EARL, JENNIFER, ANDREW MARTIN, JOHN D. MCCARTHY, AND SARAH A. SOULE.**

2004 The Use of Newspaper Data in the Study of Collective Action. *Annual Review of Sociology* 30:65-80.

**ECK, KRISTINE.**

2012 In Data We Trust? A Comparison of UCDP GED and ACLED Conflict Event Datasets. *Cooperation and Conflict* 47:124-41.

**EISINGER, PETER K.**

1973 The Conditions of Protest Behavior in American Cities. *The American Political Science Review* 67:11-28.

**ENTMAN, ROBERT M.**

1993 Framing: Toward Clarification of a Fractured Paradigm. *Journal of Communication* 43:51-58.

**FEARON, JAMES D. AND DAVID D. LAITIN.**

2003 Ethnicity, Insurgency, and Civil War. *American Political Science Review* 97:75-90.

**FRANZOSI, ROBERTO.**

1987 The Press as a Source of Socio-historical Data: Issues in the Methodology of Data Collection from Newspapers. *Historical Methods: A Journal of Quantitative and Interdisciplinary History* 20:5-16.

**GALTUNG, JOHAN AND MARI H. RUGE.**

1965. The Structure of Foreign News: The Presentation of the Congo, Cuba and Cyprus Crises in *Four Norwegian Newspapers*. *Journal of Peace Research* 2:64-90.

**GROVES, ROBERT M.**

1989. *Survey Errors and Survey Costs*. New York: Wiley.

**GROVES, ROBERT M., FLOYD J. FOWLER, MICK P. COUPER, JAMES M. LEPKOWSKI, ELEANOR SINGER, AND ROGER TOURANGEAU.**

2004 *Survey Methodology*. Hoboken, NJ: Wiley.



**GLEDITSCH, KRISTIAN S., NILS W. METTERNICH, AND ANDREA RUGGERI.** 2014 Data and Progress in Peace and Conflict Research. *Journal of Peace Research* 51:301-14.

**HAMMOND, JESSE AND NILS B. WEIDMANN.**  
2014 Using Machine-coded Event Data for the Micro-level Study of Political Violence. *Research & Politics* 1:1-8.

**HARCUP, TONY AND DEIRDRE O'NEILL.**  
2001 What Is News? Galtung and Ruge Revisited. *Journalism Studies* 2:261-80.

**HARCUP, TONY AND DEIRDRE O'NEILL.**  
2016 What Is News? News Values Revisited (Again). *Journalism Studies* 18:1-19. doi: 10.1080/1461670X.2016.1150193.

**HEAP, BRADFORD, ALFRED KRZYWICKI, SUSANNE SCHMEIDL, WAYNE OBCKE, AND MICHAEL Bain**  
2017. "A Joint Human/Machine Process for Coding Events and Conflict Drivers." Conference Paper, International Conference on Advanced Data Mining and Applications, ADMA 2017, Singapore. October.

**HENDRIX, CULLEN S. AND IDEAN SALEHYAN.**  
2015. No News Is Good News: Mark and Recapture for Event Data When Reporting Probabilities Are Less Than One. *International Interactions* 41:392-406.

**HENDRIX, CULLEN S. AND IDEAN SALEHYAN.**  
2017. A House Divided: Threat Perception, Military Factionalism, and Repression in Africa. *Journal of Conflict Resolution*. 61:1653-81.

**HERKENRATH, MARK AND ALEX KNOLL.**  
2011. Protest Events in International Press Coverage: An Empirical Critique of cross-national Conflict Databases. *International Journal of Comparative Sociology* 52:163-80.

**HICKLER, CHRISTIAN AND JOHN WIESEL.**  
2012. New Method, Different War? Evaluating Supervised Machine Learning by Coding Armed Conflict. *SFB-Governance Working Paper Series*, No.39. Collaborative Research Center (SFB) 700, Berlin. 2012. Retrieved October 12, 201

- Muestreo de eventos relevantes<sup>22</sup>: 80-97% de recuperación y 23-58% de precisión.  
- Categorización correcta de los eventos<sup>23</sup>: 7-96 por ciento.  
(<https://refubium.fuberlin.de/bitstream/handle/fub188/18738/WP39.pdf?sequence=41&isAllowed=1>).

**HIRSHLEIFER, JACK.**

<sup>22</sup> Bond et al. (1997), Croicu y Weidmann (2016), King y Lowe (2003).

<sup>23</sup> Bond et al. (1997), Boschee et al. (2013), King y Lowe (2003), Stepinski et al. (2006).

Demarest, Leila y Langer, Arnim "Cómo se incorporan (o no) los eventos en los sets de datos: Los obstáculos y las orientaciones en la utilización de los periódicos en el estudio de los conflictos" *Revista de Estudios Marítimos y Sociales*, N°24, ene-jun 2024, pp. 197-244.



1994. The Dark Side of the Force. *Economic Inquiry* 32:1-10.

**HUG, SIMON AND DOMINIQUE WISLER.**

1998. Correcting for Selection Bias in Social Movement Research. *Mobilization: An International Quarterly* 3:141-61.

**HUTTER, SVEN.**

2014. Protest Event Analysis and Its Offspring. in *Methodological Practices in Social Movement Research*, edited by Donatella della Porta. Oxford, UK: Oxford University Press. 335-67

**JENKINS, J. CRAIG AND THOMAS V. MAHER.**

2016. What Should We Do about Source Selection in Event Data? Challenges, Progress, and Possible Solutions. *International Journal of Sociology* 46:42-57.

**KALYVAS, STATHIS N.**

2006. The Logic of Violence in Civil War. Cambridge, MA: Cambridge University Press.

**KING, GARY AND WILL LOWE.**

2003. An Automated Information Extraction Tool for International Conflict Data with Performance as Good as Human Coders: A Rare Events Evaluation Design. *International Organization* 57:617-42.

**KOOPMANS, RUUD AND DIETER RUCHT.**

2002. Protest Event Analysis. in *Methods of Social Movement Research*, edited by B. Klandermans and S. Staggenborg. Minneapolis: University of Minnesota 31-59

**KRIESI, HANSPETER, RUUD KOOPMANS, JAN W. DUYVENDAK, AND MARCO G. GIUGNI.**

1998. New Social Movements in Western Europe: A Comparative Analysis. Minneapolis: University of Minnesota.

**KRIPPENDORFF, KLAUS.**

2013. Content Analysis: An Introduction to Its Methodology. *Thousand Oaks, CA: Sage.*

**LEETARU, KALEV AND PHILIP A. SCHRODT.**

2013. GDELT: Global Data on Events, Location and Tone, 1979-2012. *Paper presented at the International Studies Association Meetings*, San Francisco, CA, April.

**LAFREE, GARY AND LAURA DUGAN.**

2007. Introducing the Global Terrorism Database. *Terrorism and Political Violence* 19:181-204.

**LEE, FRANCIS L. F.**

2014. Triggering the Protest Paradigm: Examining Factors Affecting News Coverage of Protests. *International Journal of Communication* 8: 2725-46.

**MCCARTHY, JOHN D., CLARK MCPHAIL, JACKIE SMITH, AND LOUIS J. CRISHOCK.**

Demarest, Leila y Langer, Arnim "Cómo se incorporan (o no) los eventos en los sets de datos: Los obstáculos y las orientaciones en la utilización de los periódicos en el estudio de los conflictos" *Revista de Estudios Marítimos y Sociales*, N°24, ene-jun 2024, pp. 197-244.





1999. Electronic and Print Media Representations of Washington D.C. Demonstrations, 1982 and 1991: A Demography of Description Bias in *Acts of Dissent: New Developments in the Study of Protest*, edited by Dieter Rucht, Ruud Koopmans, and Friedhelm Neidhart. Oxford, England: Rowman & Littlefield. 113-30

**MCCLELLAND, CHARLES A.**

1976. World Event/Interaction Survey Codebook (ICPSR 5211). Ann Arbor, MI: Inter-University Consortium for Political and Social Research.

**MCLEOD, DOUGLAS AND JAMES HERTOGL.**

1992. The Manufacture of 'Public Opinion' by Reporters: Informal Cues for Public Perceptions of Protest Groups. *Discourse & Society* 3:259-75.

**MYERS, DANIEL J. AND BETH S. CANIGLIA.**

2004. All the Rioting That's Fit to Print: Selection Effects in National Newspaper Coverage of Civil Disorders, 1968-1969. *American Sociological Review* 69:519-43.

**OLIVER, PAMELA E. AND DANIEL J. MYERS.**

1999. How Events Enter the Public Sphere: Conflict, Location, and Sponsorship in Local Newspaper Coverage of Public Events. *American Journal of Sociology* 105:38-87.

**ORTIZ, DAVID G., DANIEL J. MYERS, N. EUGENE WALLS, AND MARIA-ELENA D. DIAZ.**

2005. Where Do We Stand with Newspaper Data? *Mobilization: An International Journal* 10:397-419.

**RALEIGH, CLIONADH AND CAITRIONA DOWD.**

2017. Armed Conflict and Event Location (ACLED) Codebook. Retrieved 12 October 2019  
([https://www.acleddata.com/wpcontent/uploads/2017/01/ACLED\\_Codebook\\_2017.pdf](https://www.acleddata.com/wpcontent/uploads/2017/01/ACLED_Codebook_2017.pdf))

**RALEIGH, CLIONADH, ANDREW LINKE, HARVARD HEGRE, AND JOAKIM KARLSEN.**

2010. Introducing ACLED-Armed Conflict Location and Event Data. *Journal of Peace Research* 47:1-10.

**REINEMANN, CARSTEN, JAMES STANYER, SEBASTIAN SCHERR, AND GUIDO LEGNANTE.**

2011. Hard and Soft News: A Review of Concepts, operationalizations and Key Findings. *Journalism* 13:221-39.

**RUGGERI, ANDREA, THEODORA-ISMENE GIZELIS, AND HAN DORUSSEN.**

2011. "Events Data as Bismarck's Sausages? Intercoder Reliability, Coders' Selection, and Data Quality." *International Interactions* 37:340-61.

**SALEHYAN, IDEAN.**

Demarest, Leila y Langer, Arnim "Cómo se incorporan (o no) los eventos en los sets de datos: Los obstáculos y las orientaciones en la utilización de los periódicos en el estudio de los conflictos" *Revista de Estudios Marítimos y Sociales*, N°24, ene-jun 2024, pp. 197-244.



2015. Best Practices in the Collection of Conflict Data. *Journal of Peace Research* 52:105-9.

**SALEHYAN, IDEAN, CULLEN S. HENDRIX, JESSE HAMNER, CHRISTINA CASE, CHRISTOPHER LINEBARGER, EMILY STULL, AND JENNIFER WILLIAMS.**

2012. Social Conflict in Africa: A New Database. *International Interactions* 38:503-11.

**SCHRODT, PHILIP A.**

2006. Twenty Years of the Kansas Event Data System Project. *The Political Methodologist* 14:2-8.

**SCHRODT, PHILIP A.**

2012. Precedents, Progress, and Prospects in Political Event Data. *International Interactions* 38:546-69.

**SCHRODT, PHILIP A. AND DAVID VAN BRACKLE.**

2013. Automated Coding of Political Event Data. in *Handbook of Computational Approaches to Counterterrorism*, edited by V. S. Subrahmanian. New York: Springer Science & Business Media. 23-49

**SCHRODT, PHILIP A., ERIN M. SIMPSON, AND DEBORAH J. GERNER.**

2001. Monitoring Conflict Using Automated Coding of Newswire Reports: A Comparison of Five Geographical Regions. Conference paper, June 8-9, Uppsala, Sweden.

**SCHRODT, PHILIP A. AND JAY ULFELDER.**

2016. Political Instability Task Force Atrocities Event Data. *Collection Codebook Version 1.1b1*. Retrieved October 12, 2019([http://eventdata.parusanalytics.com/data.dir/PITF\\_Atrocities.codebook.1.1B1.pdf](http://eventdata.parusanalytics.com/data.dir/PITF_Atrocities.codebook.1.1B1.pdf))

**SMITH, JACKIE, JOHN D. MCCARTHY, CLARK MCPHAIL, AND AUGUSTYN BOGUSLAW.**

2001. From Protest to Agenda Building: Description Bias in Media Coverage of Protest Events in Washington, D.C. *Social Forces* 79:1397-423.

**SCHNEIDER, GERALD AND MARGIT BUSSMANN.**

2013. Accounting for the dynamics of onesided violence: Introducing KOSVED. *Journal of Peace Research* 50:635-44.

**STEPINSKI, ADAM, RICHARD STOLL, AND DEVIKA SUBRAMANIAN.**

2006. Automated Event Coding Using Conditional Random Fields. Retrieved June 13, 2018 (<https://www.cs.rice.edu/~devika/conflict/papers/draft1.pdf>).

**SUNDBERG, RALPH AND ERIK MELANDER.**

2013. Introducing the UCDP Georeferenced Event Dataset. *Journal of Peace Research* 50:523-32.

**URDAL, HENRIK.**



2008. Urban Social Disturbance in Africa and Asia: report on a New Dataset. *PRIO Papers*, Oslo, Norway.

**WEAVER, DAVID A. AND JOSHUA M. SCACCO.**

2012. Revisiting the Protest Paradigm. *The International Journal of Press/Politics* 18:61-84.

**WEIDMANN, NILS B.**

2015. On the Accuracy of Media-based Conflict Event Data. *The Journal of Conflict Resolution* 59:1129-49.

**WEIDMANN, NILS B.**

2016. A Closer Look at Reporting Bias in Conflict Event Data. *American Journal of Political Science* 60:206-18.

**WEIDMANN, NILS B. AND ESPEN G. RØD.**

2015. Making Uncertainty Explicit: Separating Reports and Events in the Coding of Violence and Contention. *Journal of Peace Research* 52:125-28.

**WILKINSON, STEVEN I.**

2004. Votes and Violence: Ethnic Competition and Ethnic Riots in India. New York: Cambridge University Press.

### Biografías de los autores

Leila Demarest es profesora adjunta de Política Africana en el Departamento de Ciencias Políticas de la Universidad de Leiden (Países Bajos). Sus intereses de investigación incluyen los movimientos sociales y la movilización política en África, la comunicación política y la metodología cuantitativa y cualitativa de la investigación en ciencias sociales. Sus publicaciones más recientes (con Arnim Langer) son "The Study of Violence and Social Unrest in Africa: A Comparative Analysis of Three Conflict Event Datasets" en *African Affairs* (abril de 2018) y "Peace Journalism on a Shoestring? Conflict Reporting in Nigeria's National News Media" en *Journalism* (primero en línea, agosto de 2018).

Arnim Langer es profesor de Relaciones Internacionales en la Universidad Católica de Lovaina y director del Centro de Investigación sobre la Paz y el Desarrollo (CRPD) de la Facultad de Ciencias Sociales. Actualmente, es también investigador Humboldt en la Universidad de Heidelberg, Alemania. Ha publicado numerosos trabajos sobre las causas de los conflictos violentos en sociedades heterogéneas y los retos que plantea la consolidación de una paz sostenible. Algunas de sus publicaciones más recientes son

Demarest, Leila y Langer, Arnim "Cómo se incorporan (o no) los eventos en los sets de datos: Los obstáculos y las orientaciones en la utilización de los periódicos en el estudio de los conflictos" *Revista de Estudios Marítimos y Sociales*, N°24, ene-jun 2024, pp. 197-244.





"Conceptualising and Measuring Social Cohesion in Africa: Towards a Perceptions-based Index" (publicado en Social Indicators Research) y "A General Class of Social Distance Measures" (publicado en Political Analysis).