

dosr: Acelerando y Estandarizando el Análisis de Encuestas Complejas en el Sector Público

Palabras clave: srvyr, encuestas complejas, sector público, automatización, reproducibilidad.

Gabriel Sotomayor López

Abstract

Las instituciones del sector público dependen del análisis de encuestas con diseños muestrales complejos para informar políticas públicas. Este proceso suele ser manual, con posibilidad de errores y difícil de estandarizar. Para abordar estos desafíos, dentro de la División de Observatorio Social (DOS) del Ministerio de Desarrollo Social y Familia de Chile se desarrolló dosr, un paquete de R que encapsula el flujo de trabajo completo. Sus funciones de alto nivel (obs_prop, obs_media, etc.), integradas con el ecosistema srvyr, automatizan el cálculo de estimaciones para múltiples encuestas simultáneamente, siendo aplicable a cualquier conjunto de encuestas complejas que el usuario prepare como objetos tbl_svy. Sus características principales son la implementación programática de criterios de fiabilidad estadística institucionales y el cálculo automático de pruebas de significancia. Finalmente, dosr genera reportes en Excel listos para difusión. Esta charla presenta dosr como un caso de estudio sobre cómo R puede aumentar la eficiencia y estandarizar la calidad de la evidencia estadística en el servicio público.

Introducción: El Desafío del Análisis de Encuestas en el Sector Público

El análisis riguroso de encuestas nacionales es una tarea fundamental para el monitoreo social en el sector público. Este proceso requiere no solo el cálculo de estimadores, sino también la evaluación de su calidad estadística para asegurar la robustez de la evidencia. Tradicionalmente, este flujo de trabajo no está estandarizado, abriendo espacio a inconsistencias y requiriendo la aplicación manual de criterios de fiabilidad (basados en grados de libertad y errores estándar o coeficiente de variación) y la programación ad-hoc de pruebas de significancia.

Para enfrentar estos problemas, se desarrolló una solución interna utilizando R, lo cual se suma a una solución previa desarrollada en STATA, con el objetivo de automatizar, estandarizar y acelerar el proceso completo, desde el cálculo hasta el reporte final.

La Solución: El Paquete dosr

La solución implementada fue el paquete de R dosr. Esta herramienta proporciona una interfaz de alto nivel que abstrae la complejidad del análisis de encuestas, permitiendo a los analistas enfocarse en la interpretación de resultados. Integrado con srvyr, el paquete ofrece cinco funciones principales, obs_prop(), obs_media(), obs_cuantil(), obs_total y obs_ratio() que con una sola llamada ejecutan el ciclo completo de análisis.

Ejemplo: Media de ingresos por sexo y área para tres años, con tests de significancia

```
obs_media(designs = list("2017" = casen_2017d, "2020" = casen_2020d, "2022" = casen_2022d),
          var = "ytrabajocorh",
          des = c("sexo", "area", "region"),
          parallel = T,
          formato = T,
          multi_des = T,
          filt = "pco1==1 & nucleo!=0",
          sig = TRUE)
```

Características y Flujo de Trabajo

El diseño de dosr encapsula las mejores prácticas del análisis de encuestas en una herramienta simple y potente:

- Suite Completa de Estimadores: Calcula proporciones, medias, cuantiles, totales y ratios a través de una API coherente.
 - Manejo de Múltiples Encuestas: Facilita el análisis comparativo y de tendencias al aceptar una lista de objetos `tbl_svy`. Para garantizar comparaciones válidas, es responsabilidad del usuario asegurar la estandarización de variables y etiquetas entre los diferentes diseños.
 - Desagregaciones Flexibles: El usuario puede solicitar cruces por múltiples variables (`multi_des = TRUE`) o solo desagregaciones simples (`multi_des = FALSE`).
 - Implementación de Criterios de Calidad: El motor interno aplica automáticamente el flujo de decisión para la fiabilidad de las estimaciones, añadiendo una columna fiabilidad ("Fiable", "Poco Fiable", "No Fiable") a todos los resultados.
 - Pruebas de Significancia Automatizadas: Con `sig = TRUE`, el paquete calcula p-values para tres tipos de comparaciones: entre categorías, contra el total nacional y a través del tiempo.
 - Generación de Reportes: La salida principal es un archivo Excel con una hoja de datos consolidados y múltiples hojas con tablas de presentación formateadas, listas para su difusión (ver Figura 1).

A	B	C	D	E	F	G	H	I	J	K	L	
variable	nivel	area	region	media_2012	se_2012	cv_2012	b_mues_2011	n_pob_2017	pl_2017	abilidad_2017	...	
2	trabajador	Nacional		94115,824	14330,4131	0,01703649	79940	599774	1295	Flujo		
3	trabajador	sexo	1. Hombre	93667,754	17484,6686	0,01826213	41450	3670162	1295	Flujo		
4	trabajador	sexo	2. Mujer	65948,283	13651,7259	0,02070017	29498	237570	1295	Flujo		
5	trabajador	area	Urbano	97304,754	14764,3646	0,01826213	57940	5294200	1295	Flujo		
6	trabajador	area	Rural	50480,871	8000,4779	0,02007009	13485	12500	1294	Flujo		
7	trabajador	region	Región de D	31013,276	400800011	2953	10148	56	Flujo			
8	trabajador	region	Región de T	10470,54	45966,534	0,02842585	2628	202549	49	Flujo		
9	trabajador	region	Región de A	789284,378	46234,9311	0,0587829	2348	97767	40	Flujo		
10	trabajador	region	Región de C	57409,397	57529,325	0,01010519	3127	245377	56	Flujo		
11	trabajador	region	Región de E	59052,796	57035,6398	0,02036181	5079	412790	86	Flujo		
12	trabajador	region	Región de L	59352,796	57035,6398	0,02036181	5244	312790	88	Flujo		
13	trabajador	region	Región del N	57684,498	19448,251	0,0336241	5140	361499	86	Flujo		
14	trabajador	region	Región del B	61697,844	37349,5406	0,06062529	7177	54757	135	Flujo		
15	trabajador	region	Región de M	57594,829	25813,019	0,0448613	5189	331246	91	Flujo		
16	trabajador	region	Región de S	50364,363	34617,4627	0,051224	4160	28000	73	Flujo		
17	trabajador	region	Región de R	57650,361	34617,4627	0,051224	3779	32744	36	Flujo		
18	trabajador	region	Región de P	102762,799	80580,789	0,09520572	3232	59572	46	Flujo		
19	trabajador	region	Región de M	101483,83	37208,0827	0,02287699	24030	2405660	247	Flujo		
20	trabajador	region	Región de G	631548,756	51607,5996	0,08125278	3403	132793	67	Flujo		
21	trabajador	region	Región de A	626314,878	30551,7245	0,05387907	2617	17799	53	Flujo		
22	trabajador	region	Región de C	626314,878	30551,7245	0,05387907	2617	17799	53	Flujo		
23	trabajador	sex-area	1. Hombre	100725,46	20000,46	0,01811315	3203	316170	1298	Flujo		
24	trabajador	sex-area	1. Hombre	635857,599	55398,3706	0,01821298	9245	503044	213	Flujo		
25	trabajador	sex-area	2. Mujer	678600,906	14311,3973	0,02108061	25255	2127180	1082	Flujo		
26	trabajador	sex-area	2. Mujer	456455,309	4471,759	0,09729559	424	20040	213	Flujo		
1. Consolidado												
2. Ingreso del trabajo del hogar												
3. Tipo de cálculo												
4. Estimación												
5. Test entre categorías año: 2017												
area			2017			2020			2022			
Urbanos			875.221,52			849.073,13			1.122.640,16			
Rurales			584.000,58			575.889,99			575.271,92			
Total país			1.419.320,50			1.424.961,12			1.598.608,08			
6. Test entre categorías año: 2020												
area			2017			2020			2022			
Urbanos			529.351,76			510.351			494.975,73			
Rurales			459.944,74			458.914			455.118,15			
Total país			999.742			968.673,27			998.093			
7. Test entre categorías año: 2022												
area			2017			2020			2022			
Urbanos			1.416.746,30			1.405.197,29			1.288.161,16			
Rurales			587.813,71			575.457,49			575.873,89			
Total país			1.404.330,11			1.400.653,11			1.176.708,64			
8. Error estándar												
area			2017			2020			2022			
Urbanos			14.764			14.029			13.286,16			
Rurales			587.813,71			575.457,49			575.873,89			
Total país			1.419.320,50			1.400.653,11			1.176.708,64			
9. Test contra estimación final:												
sexo			area			region			p_value_2017			
Urbanos			Rural			Rural			0,223			
Rurales			Urbanos			Urbanos			0,132			
10. Test contra estimación final:												
area			2017			2020			2022			
Urbanos			57.460			52.993			51.130			
Rurales			13.488			9.918			14.928			
Total país			70.948			62.911			72.056			

Figura 1: Ejemplo de hoja de reporte generada por dosr, mostrando estimaciones y pruebas de significancia.

Impacto y Desafíos Futuros

La adopción de dosr, junto con otras herramientas desarrolladas por el equipo, ha tenido un impacto directo en la eficiencia y calidad del trabajo de análisis. Las tareas se realizan de manera más eficiente y estandarizada, entregando una mejor herramienta para asegurar que todas las estimaciones publicadas son sometidas a los mismos estándares de fiabilidad. Esto no solo mejora la reproducibilidad, sino que también facilita la incorporación de nuevos analistas al equipo.

El paquete dosr es un ejemplo de cómo el desarrollo de herramientas de código abierto dentro del sector público puede mejorar sustancialmente los flujos de trabajo. Los desafíos y planes futuros se centran en:

- Implementar una suite de pruebas unitarias (testthat) para garantizar la estabilidad a largo plazo.
 - Generar un sitio web del paquete (pkgdown) para mejorar la documentación y facilitar su adopción.
 - Abstraer los criterios de fiabilidad para que puedan ser personalizados por el usuario, aumentando la utilidad del paquete en otros contextos institucionales.