

Seminario interno UNA - 17/11/2025

Gestión de proyectos de investigación reproducibles:

Federico Giovannetti
giovannettipsi@gmail.com



Unidad de Neurobiología Aplicada



CEMIC-CONICET

Motivaciones

- Crisis de replicabilidad/reproducibilidad
- Necesidad de construir una ciencia más abierta y transparente
- Necesidad de mejorar flujos de trabajo a nivel individual y colectivo.

Crisis de replicabilidad

Pensemos en alguno de nuestros trabajos de investigación...

Si alguien quisiera repetir nuestro estudio siguiendo la misma metodología pero con nuevos datos similares (aparentemente a los nuestros) ¿Obtendrían los mismos resultados?

Crisis de replicabilidad

Para que pueda haber replicabilidad, debería haber algún nivel de *invarianza*.

Sin embargo, en psicología y ciencias afines sabemos que la invarianza es difícil (¿imposible?) de conseguir.

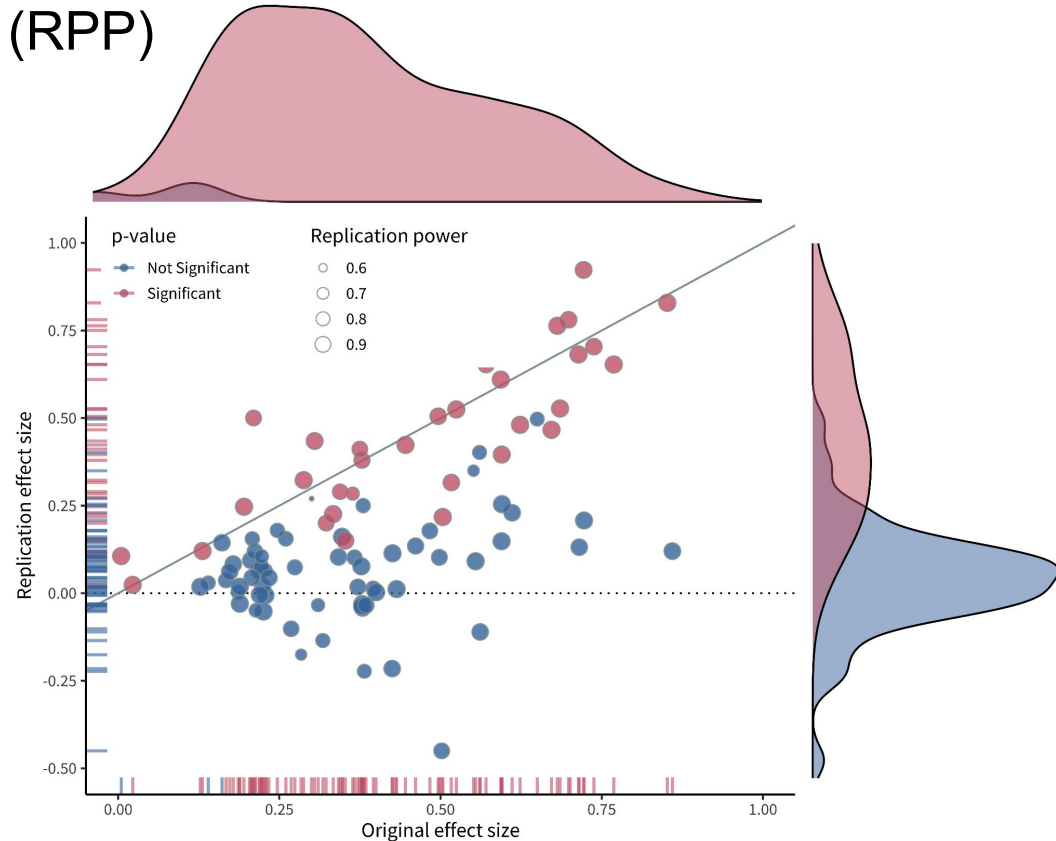
Crisis de replicabilidad

Replication Project in Psychology (RPP)

En 2015 se publicó un trabajo donde distintos grupos de estudiantes intentaron replicar 100 estudios psicológicos.

Solo un tercio de los trabajos pudieron ser replicados.

En general, hubo tamaños del efecto más pequeños de lo esperado.



Crisis de replicabilidad

Existen un montón de discusiones interesantes sobre el tema de la replicabilidad:

Sus causas

Sus consecuencias

Sus implicancias teóricas

Sobre si cuán bien nos cae Popper

etc

Pero lo que me interesa profundizar en este
seminario es en la parte más práctica

Reproducibilidad

La investigación científica implica **mucho** trabajo

En general, “termina” con algún tipo de publicación o presentación donde se muestra una serie de resultados, números, gráficos, etc.

Sin embargo:

- ¿Cuánto del proceso de investigación logra ser reflejado?
- ¿Alcanza con la descripción **narrativa** de diseño, participantes, instrumentos, análisis, etc.?
- ¿Alcanza la información disponible para que otro grupo reproduzca los resultados?
- ¿Qué otra información debería haber para que el trabajo sea efectivamente reproducible?

Reproducibilidad



Incluso teniendo los datos

En un trabajo de Hardwicke, Bohn, y colegas, (2021), intentaron reproducir los valores numéricos reportados en una serie de artículos con datos abiertos.

Sólo un tercio de los artículos fueron completamente reproducibles sin la ayuda de los autores

En otros casos, se precisó una comunicación intensiva con lxs autores para llegar al resultado final

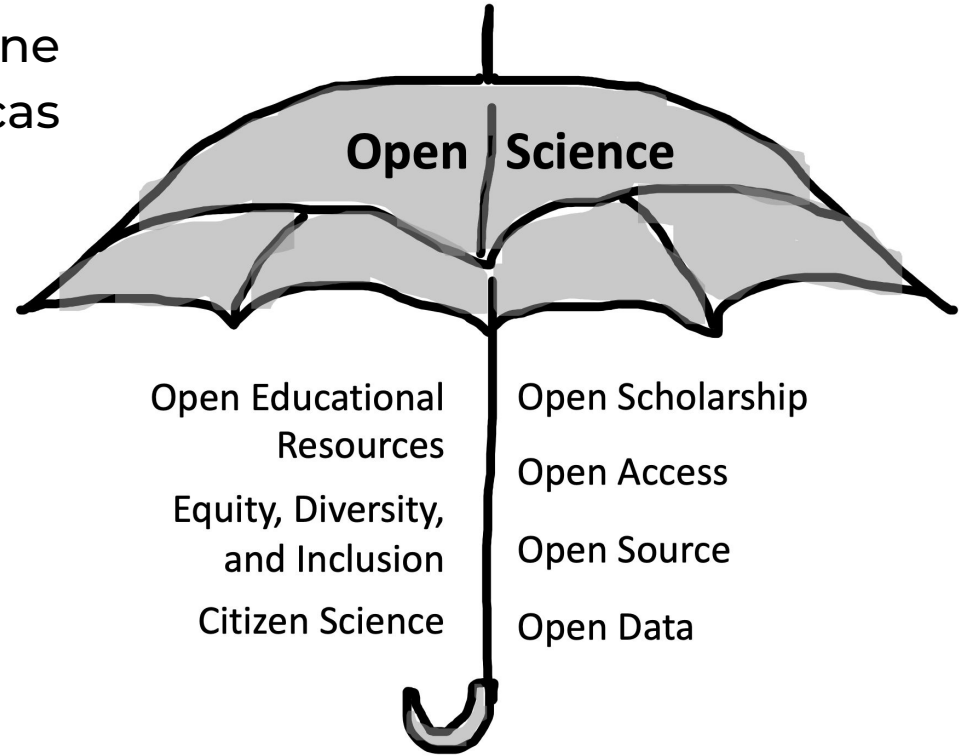
Otros trabajos no pudieron ser reproducidos ni por sus propios autores.

Ciencia abierta

Sobran los motivos para ser reproducibles (?)

El movimiento de ciencia abierta tiene que ver con ideas, prácticas, y políticas orientadas a la transparencia, la verificabilidad y la democratización del conocimiento y de la práctica científica.

Para todo esto, la **transparencia** y la **accesibilidad** es clave.



Ciencia abierta

Hay que compartir

Ser transparente y compartir de forma abierta los materiales de un trabajo científico podría permitir

- La replicación del trabajo realizado por otra gente (o uno mismo)
- Disminuir las probabilidades de que alguien ponga resultados truchos o plantee hipótesis a posterior (e.g. con preregistro)
- **Reducir errores en el flujo de trabajo**
si estás haciendo todo con la idea de que lo va a leer otra persona, probablemente seas más cuidadoso, más claro, y más amigable con tu yo del futuro
- **Trabajar menos**
(e.g. realizando los análisis en el mismo formato en que luego los vas a publicar)
https://github.com/FedeGiovannetti/Poster_AACC_2025

Vamos con algunas propuestas concretas

1. Flujos de trabajo reproducibles

Organizar individual y colectivamente el flujo de trabajo de forma explícita puede ordenar el trabajo y las expectativas.

Your closest collaborator is you six months ago, but you don't reply to emails.

—Karl Broman (2015), quoting @gonuke on Twitter

1. Flujos de trabajo reproducibles

¿Cuáles son los distintos momentos por los que pasa un proyecto de investigación?

1. Flujos de trabajo reproducibles

Diseño instrumentos

- ¿Cuáles son los cuestionarios a utilizar?
- ¿Qué tareas cognitivas?
- ¿Qué estímulos?
- ¿Qué softwares?
- ¿Qué equipos?

Procedimiento

- ¿Quiénes van a implementar el trabajo?
- ¿Qué capacitación recibieron?
- ¿Hay algún manual de procedimientos?

Toma de datos

- ¿Cómo se almacenarán los datos?
- ¿La información almacenada es (al menos) la requerida?
- ¿Qué nombre tendrán los distintos archivos? ¿Cómo se llamarán las variables? ¿Cómo se anonimiza a los sujetos?

Análisis de resultados

- ¿Cuáles serán las distintas etapas por las que pasarán los datos?
- ¿Cómo se llama cada variable? ¿Cómo se codifican? ¿Hay un codebook?
- ¿Con qué programa se realizará el análisis?
- ¿Quién(es) van a encargarse de los análisis? ¿Habrá revisión del código?

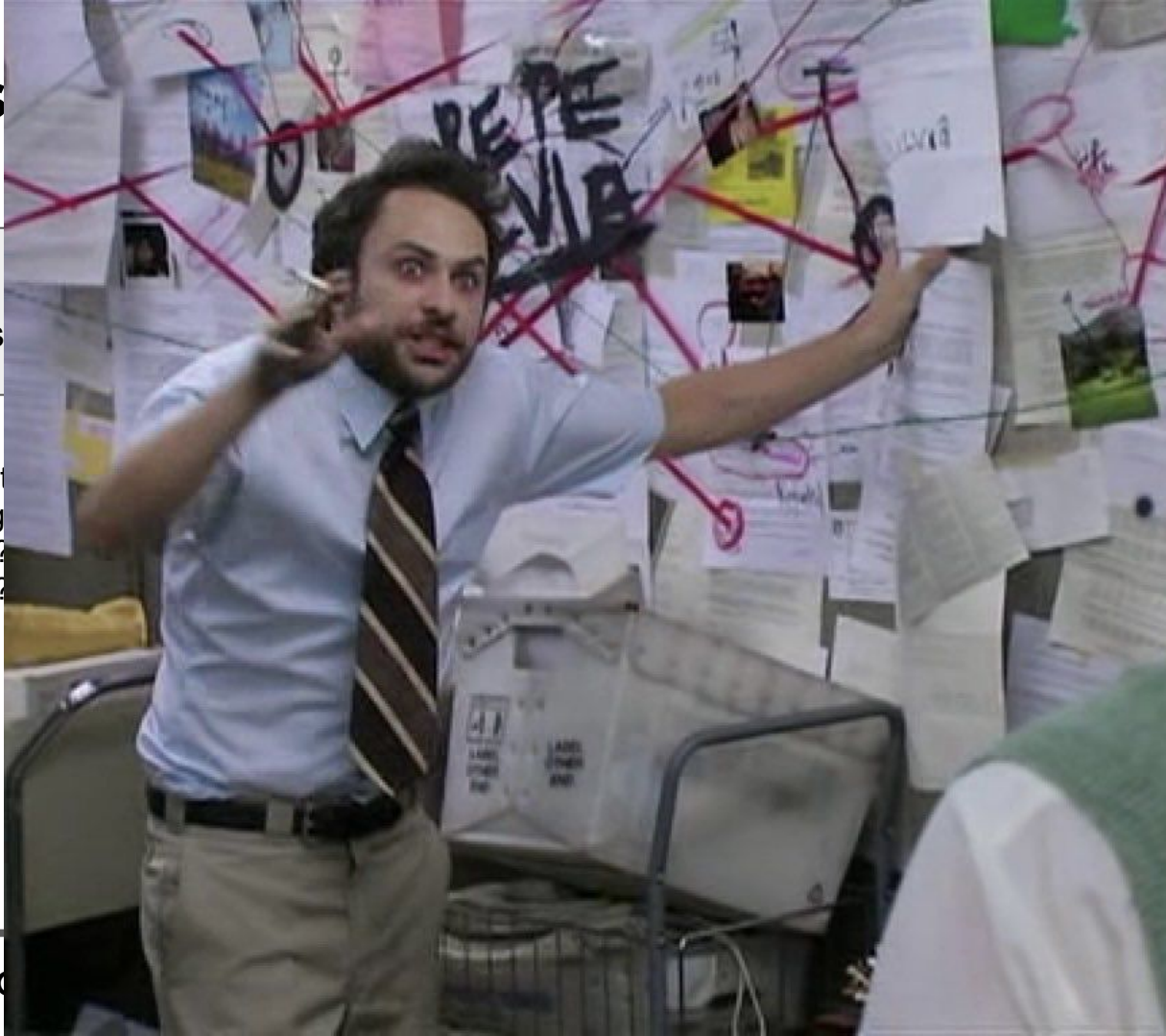
¿Dónde va a estar registrado todo esto? ¿De qué forma?

1. Flujos

Diseño
instrumentos

- ¿Cuáles son los cuestionarios a utilizar?
- ¿Qué tareas cognitivas?
- ¿Qué estímulos?
- ¿Qué softwares?
- ¿Qué equipos?

¿De



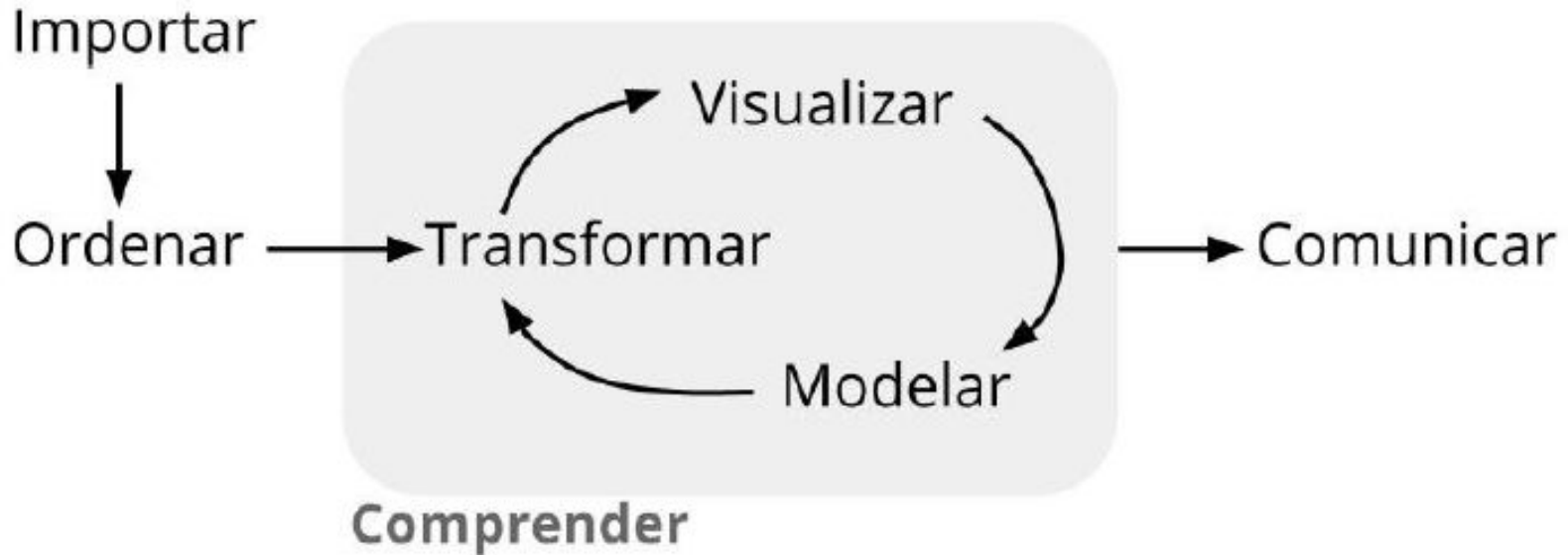
Análisis de
resultados

¿Cuáles serán las distintas tareas por las que pasarán los usuarios?
¿Cómo se llama cada uno de ellos?
¿Cómo se llaman los datos?
¿Hay un libro de usuario?
¿Qué programa se utilizará para el análisis?
¿Quién(es) van a realizar el análisis?
¿Se hará una revisión del código?

¿De

1. Flujos de trabajo reproducibles

Para análisis de datos, se podría pensar algo así:



2. Documentar todo

En general, terminamos teniendo archivos tipo: `datos_1.csv`, `datos_2.csv`, `datos_completos.csv`, `datos_completos_ahora_si.xlsx`, etc....

Para ordenar todo esto, hay que **documentar** nuestro flujo de trabajo y generar archivos de **metadatos** (es decir, de datos sobre nuestros datos).

Y cuando hablamos de datos nos referimos no solo a las tablas, figuras, etc. También entran los materiales (e.g. tareas, estímulos, manuales) que se usaron en todo el flujo de trabajo

Formatos comunes son los **ReadMe**, los **codebooks** (diccionarios de código), pero básicamente son documentos que expliquen **dónde están** las cosas, **qué son** cada cosa, **cómo se llaman** las cosas, **en qué formato están**, etc.

2. Documentar todo

Un primer intento...

Tabla_1											
Tarea	Tiempo estimado	Inicio	Deadline	Estado	Responsable	Revisor	Sub-pipeline	Materiales			
Adaptación de los cuentos			1/9/2025	Terminado	Fede	Fede Marcos	Escucha dicótic...	Cuentos			
Grabación de los cuentos	Gada cuento: 30 mins. 18 h	01/09/2025	22/09/2025	Cancelado	lae	Fede Marcos		audios_viejos			
Edición de los audios	Charlar con Juan qué había	22/09/2025	13/10/2025	Cancelado	Fede	Marcos					
Generación de audios con IA				En proceso	lae	Fede Marcos	Escucha dicótic...	audios_ElevenLa...			
Preparación del paradigma				En proceso	Juan Octavio	Marcos Fede					
Prueba piloto				Pendiente	lae	Fede Marcos					
Preparación de la toma de datos				Pendiente	lae	Fede					
Proyecto de tesis (2da fecha)			1ra semana marzo	Pendiente	lae	lae Fede					
Vacaciones											
Toma de datos		1/2/2025	1/4/2026	Pendiente	lae	Fede Marcos					
Procesamiento de datos	Va en simultaneo con la tom		01/05/2026	Pendiente							
Análisis de datos			01/07/2025	Pendiente							
Tesis			Principios agosto	Pendiente	lae	Fede Marcos					
Tesis (2da fecha)			Principios noviembre	Pendiente	lae	Fede Marcos					

Esto permite **reconstruir** la cadena de acciones del proyecto

3. Trabajar con control de versiones

A veces te ponés a trabajar en un análisis, o en unos datos, o en un escrito y no sabes si estás trabajando en la versión final, la versión_final_2, etc.

GoogleDocs podría ser una forma de control de versiones ya que todo el mundo trabaja ahí y disminuye la posibilidad de que haya archivos duplicados.

Para cosas de código y análisis, sistemas como **git** permiten ir documentando los cambios que hacemos en nuestro trabajo, que los demás lo vean y que, quizás, sumen lo suyo.

4. Tener criterios consistentes para organizar los datos

A veces en el apuro y el pluriempleo vamos generando carpetas, archivos, nombres de variables y demás sin seguir un criterio claro.

- No respetando mayúsculas, y minúsculas (Ponemos “Edad” en una base de datos, “edad” en otra).
- No respetando el nombre elegido anteriormente (ponemos “EDAD_NIÑO” en otra base).
- Cambiando la forma de organizar los datos (e.g. en un proyecto tenemos dividido en carpetas “datos/”, “analisis/” y en otros metemos todo junto

4. Tener criterios consistentes para organizar los datos

Eventualmente deberíamos generar criterios comunes.

Experimentology es una buena base, también los criterios *tidy* de Hadley Wickham.

A continuación, recopilo algunos aspectos relevantes:

4.1. Darle una organización jerárquica al proyecto

- Que todos los archivos estén en una misma carpeta
- Que cada sub-carpeta corresponda a distintos momentos del proyecto.

-

Name ^ v	Modified ^ v
Example project (/rpydu/)	
- OSF Storage (United States)	
+ Analyses	
Heycke, Aust, & Stahl (2017) Subliminal influence on prefer...	2018-01-12 06:29 AM
+ Material	
+ Processed data	
+ Raw data	
README.md	2018-06-12 07:26 AM
Study protocol (Stage-1 registered report).pdf	2018-01-12 06:33 AM

4.1. Darle una organización jerárquica al proyecto

- Que todos los archivos estén en una misma carpeta
- Que cada sub-carpeta corresponda a distintos momentos del proyecto.
- Que la ubicación de cada archivo sea informativa del mismo por contexto (e.g. **“datos/pre_intervención/stroop_limpio.csv”** se va a referir claramente a los datos limpios de la tarea stroop en la fase pre intervención del estudio.

4.2. Que las tablas sigan el formato *tidy*

country	year	cases	population
Afghanistan	1999	745	19987071
Afghanistan	2000	666	20595360
Brazil	1999	37737	172006362
Brazil	2000	80488	174504898
China	1999	212258	1272915272
China	2000	216766	128042583

variables

country	year	cases	population
Afghanistan	1999	745	19987071
Afghanistan	2000	666	20595360
Brazil	1999	37737	172006362
Brazil	2000	80488	174504898
China	1999	212258	1272915272
China	2000	216766	128042583

observations

country	year	cases	population
Afghanistan	1999	745	19987071
Afghanistan	2000	666	20595360
Brazil	1999	37737	172006362
Brazil	2000	80488	174504898
China	1999	212258	1272915272
China	2000	216766	128042583

values

En general, todo el estilo tidy está bueno: <https://style.tidyverse.org/>

4.3. Cada columna puede tener un solo tipo de dato

- Las columnas, siguiendo la lógica tidy, no pueden dar más de una información.

Un ejemplo sería una variable llamada “Matemática” en una base de datos de notas escolares donde una maestra pone “7”, otra pone “aprobado”. No nos da la información y, además, no vamos a poder hacer cuentas.

Lo correcto sería que haya una columna que se llame “matematica_nota”, “matematica_condicion”.

4.4. La decoración no es dato

- La mayoría de los análisis hoy son llevados a cabo con programas que no miran los colores de las celdas.
- Así que, si bien no está “mal” pintar las celdas de colores para indicar algo o facilitar la lectura, esa info también tiene que estar codificada en una columna.

¿Qué otros criterios se les ocurren?
¿Qué otros problemas habría que considerar?

¿Qué otros criterios se les ocurren?
¿Qué otros problemas habría que considerar?

¡Gracias!