

Datu zientifikoaren argitalpen eraginkorra FAIR printzipioen bidez

2021-ko Urtarrila

1 Datu zientifikoaren eskuragarritasuna kolokan

Zientziaren garapen egokia experimentuen erreprodukzioan datza: alegia, edozeinek beste ikerlariek argitaratutako experimentuak erreproduzitzeko gaitasuna edukitzea. Experimentua erreproduzitzeko pausoz-pauso berregin behar da, jatorrizko experimentua baieztatu edo baliogabetzeko. Beraz, experimentu zientifiko bat ezin bada erreproduzitu, ez dauka balio zientifikorik.

Gaur egun komunitate zientifikoko ikerlariek garatzem dituzten experimentu gehienek datuen prozesamendua dute abiapuntu eta helmuga. Horregatik, datu zientifikoak modu egokian eskuragarri egotea berebiziko garrantzia du, prozesu zientifikoak aurrera egin dezan. Tamalez, gaur egun datu zientifikoaren bererabilgarritasuna kolokan dago [1], eta beste faktore batzuekin batera, eskuragarritasun egokiaren faltak “erreprodukzio krisia” sortu du [2]: experimentu gehienak ezin dira erreproduzitu, komunitate zientifikoak ontzat emandakoak barne.

Egoera honen aurrean, ikerlari talde batek datu zientifikoak modu bererabilgarrian argitaratzeko printzipioak definitu zituen. FAIR printzipioak deitu zituzten¹: “**F**indable”, “**A**ccesible”, “**I**nteroperable”, “**R**eusable”. FAIR printzipioak lehenengo aldiz *Scientific Data* aldizkari zientifiko garrantzitsuan argitaratu ziren [3].

2 FAIR printzipioak

FAIR printzipioen laburpen ofiziala 1 irudian ikus daiteke, hurrengo azpi-ataletan printzipio bakoitza azaltzen delarik.

Kontuan hartzekoa da FAIR printzipioek ez dutela estandar bat osatzen. Alegia, printzipio **orokorrak** dira, ezin dira bete modu binarioan: ez gara inoiz helduko FAIR printzipioak osoki betetzera, beti egongo da aukera datuak modu FAIR-“ago” batean argitaratzeko. FAIR printzipioen xedea ikerlarientzako datuen argitalpenerako gida izatea da, ikerlarien erakundeek dituzten baliabideen arabera FAIR maila ezberdinak lortuz argitalpen prozesuaren bidez.

¹<https://www.go-fair.org/fair-principles/>.

Argitu beharreko beste ideia garrantzitsua FAIR printzipioek teknologiarikiko duten erlazioari dagokio. FAIR printzipioak hamarkadetan zehar baliagarriak izateko diseinatuak izan direnez, FAIR printzipioek ez dituzte teknologia zehatzak gomendatzen printzipioak gauzatzeko, nahitaez teknologiak aldakorrek baitira (Hurrengo ataletan teknologia zehatzak aipatzen direnean, implementazio adibide giza aipatzen dira, ez FAIR printzipioen parte giza). Hala ere, badago teknologia aldetik FAIR printzipioetan irizpide garrantzitsua: datuen erabilgarritasuna makinenzako erraztea dute xede (eta beraz, ondorio bezala, gizakiontzat ere). Ondorioz, FAIR printzipioek datuen kontsumitzaile nagusitzat makinak dauzkate.

Azkenik, FAIR printzipioak datu eta metadatuak dagozkiela aipatu behar da. Metadatuak (“Datuei buruzko datuak”) datuen ezaugarri orokorrak adierazten dituzte, adibidez sortze-data, egilea, datuen gaiak, eta beste hainbeste. Beraz FAIR printzipioetan “(meta)datu” hitza erabiltzen da gehienetan, “datu” eta “metadatu” hitzak banaturik soilik beharrezkoa denean aipatzen direlarik.

2.1 Aurkigarri (“Findable”)

(Meta)Datuek aurkituak izateko identifikatzaile globala behar dute, bakarra, eta ez dena gailenduko. Gaur egungo teknologiarekin, URI-ek betetzen dute funtzio hori (Uniform Resource Identifier²). Gainera, (meta)datuak argitaratzen dituen erakundeak bere gain hartu behar du sortutako identifikatzaileak mantentzeko erantzunkizuna (Ikus, adibidez, W3C erakundeak egindako promes publikoa³).

Bestalde, datuek metadatu aberatsak izan behar dituzte, kontsumitzaileek datuak aurkitzeko duten ahalmena handiagotzeko. Metadatuak datuen identifikatzailea eduki behar dute erreferentzia giza, kontsumitzailea erraz joan dadin metadatuetatik datuetara, metadatuaren bidez datuak aurkitu ondoren.

Azkenik, (meta)datuak indexatuak egon behar dute, bilatzaile orokorretan, edo zerbitzu espezializatuetan.

2.2 Eskuragarri (“Accessible”)

(Meta)Datuak eskuragarriak egon behar dute, beraien identifikatzailea erabiliz, komunikazio protokolo estandar, libre eta unibertsal baten bidez. Gaur egun HTTPS da baldintza hori betetzen duen protokolo erabiliena⁴.

Gainera, protokoloak autentifikazioa eta erabiltzaileen baimentzea bermatzeko mekanismoak behar ditu, argitaratuko diren datuek erabilpen mugak izan baititzakete (Adibidez datu klinikoak).

Azkenik, nahiz eta datuak desagertu, metadatuak eskuragarri egon behar dute, beharrezkoak ez diren bilaketak ekiditze aldera. Normalean metadatuak gordetzea datuak gordetzea baino merkeagoa eta errazagoa da, batez ere bolumen aldetik.

²<https://tools.ietf.org/html/rfc3986>.

³<https://www.w3.org/Consortium/Persistence.html>.

⁴<https://tools.ietf.org/html/rfc2818>.

To be Findable:

F1. (meta)data are assigned a globally unique and persistent identifier

F2. data are described with rich metadata (defined by R1 below)

F3. metadata clearly and explicitly include the identifier of the data it describes

F4. (meta)data are registered or indexed in a searchable resource

To be Accessible:

A1. (meta)data are retrievable by their identifier using a standardized communications protocol

A1.1 the protocol is open, free, and universally implementable

A1.2 the protocol allows for an authentication and authorization procedure, where necessary

A2. metadata are accessible, even when the data are no longer available

To be Interoperable:

I1. (meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation.

I2. (meta)data use vocabularies that follow FAIR principles

I3. (meta)data include qualified references to other (meta)data

To be Reusable:

R1. meta(data) are richly described with a plurality of accurate and relevant attributes

R1.1. (meta)data are released with a clear and accessible data usage license

R1.2. (meta)data are associated with detailed provenance

R1.3. (meta)data meet domain-relevant community standards

Irudia 1: FAIR printzipioen adierazpen ofiziala [3].

2.3 Interoperable (“Interoperable”)

(Meta)Datuek beste datuekin integratzerakoan lortzen dute baliorik handiena. Horregatik, (meta)datuak kodifikatzeko hizkuntza konputazional formala, eskuragarria, irekia, eta komunitatean erabilera zabala duena erabili behar da: adibidez, RDF (Resource Description Framework⁵) edo OWL (Web Ontology Language⁶).

Gainera, (meta)datuek FAIR printzipioak jarraitzen dituzten hiztegiekin anotatuak egon behar dute. Horrelako hiztegiak, edo ontologiak, Linked Open Vocabularies (LOV) zerbitzuan aurki daitezke⁷.

Azkenik, (meta)datuek beste (meta)datuengana jotzen duten erreferentzia kualifikatuak izan behar dituzte, beste datuekin integrazioa errazteko. Alegia, erreferentzia bezala erabiliko den estekaren bi aldean artean zer erlazio dagoen argi egon behar da (Adibidez, “parte-da”, “katalizatzen-du”, ...). Erlazioa makinenzako ulerkorra izateko aipatutako hiztegiak erabili behar dira.

2.4 Bererabilgarri (“Reusable”)

Bererabilgarriak izateko, (meta)datuak ahalik eta aberatsenak izan behar dira. Halaber, lizentzia argi eta formal bat izan behar dute (Adibidez Creative Commons⁸), eta datuen argitalpenerako aurretik existitzen diren komunitate zientifikoaren estandarrak jarraitu: adibidez, genomika funtzionaleko experimentuen datuak argitaratzeko, MINSEQE (Minimum Information about a high-throughput nucleotide SEQuencing Experiment) estandarra jarraitzea⁹.

3 FAIR printzipioak eta zientziaren etorkizuna

Sortu zirenetik FAIR printzipioen erabilpena nabarmenki zabaldu egin da: ikerketarako dirulaguntzak jasotzeko, Europar Komisioak (EC) FAIR printzipiotan oinarritutako datu kudeaketarako planak eskatzen ditu. FAIR printzipioak sektore pribatuan ere agertzen dira, hain zuzen ere Bayer¹⁰ edo Novartis¹¹ bezalako enpresa handiek datuen barne-kudeaketarako erabiltzen dituzte. Badaude ere FAIR printzipioak jarraituz datuak argitaratzen dituzten erakunde publikoak, UniProt datu-basea kasu [4].

Baina FAIR printzipioak betez datuak argitaratzea ez da erraza, azpiegitura tekniko konplexuak eraiki behar baitira [5], eta erakunde gehienek ez dituzte ez baliabide ez eta teknikari egokirik. Arazo honen aurrean badaude FAIR printzipioen arloan produktu eta zerbitzuak eskaintzen dituzten enpresak: The

⁵<https://www.w3.org/TR/rdf11-concepts/>.

⁶<https://www.w3.org/TR/owl2-overview/>.

⁷<https://lov.linkeddata.es/dataset/lov/>.

⁸<https://creativecommons.org/>.

⁹<http://fged.org/projects/minseqe/>.

¹⁰<https://www.bayer.com/>.

¹¹<https://www.novartis.com/>.

Hyve¹², FAIR Data Systems¹³, Eccenca GmbH¹⁴, ...

FAIR printzipioak gero eta garrantzi handiagoa hartuko dute, eta zientzilariek beraien datuak kudeatzeko eta argitaratzeko ahalmen handiagoa beharko dute. Hala ere, ikerlarien ahalegina emankorra izango da, datu berriak kudeatzeko eta aurkitzeko ahalmen handiagoak ikerketa hobekoak ekarriko baitituzte epe luzera.

Erreferentziak

- [1] Dominique G. Roche, Loeske E. B. Kruuk, Robert Lanfear, and Sandra A. Binning. Public data archiving in ecology and evolution: How well are we doing? *PLOS Biology*, 13(11):1–12, 11 2015.
- [2] Monya Baker. 1,500 scientists lift the lid on reproducibility. *Nature*, 533(7604):452–454, May 2016.
- [3] Mark D Wilkinson, Michel Dumontier, IJsbrand Jan Aalbersberg, Gabrielle Appleton, Myles Axton, Arie Baak, Niklas Blomberg, Jan-Willem Boiten, Luiz Bonino da Silva Santos, Philip E Bourne, et al. The FAIR Guiding Principles for scientific data management and stewardship. *Scientific data*, 3, 2016.
- [4] Leyla García Castro, Jerven Bolleman, Sebastien Gehant, Nicole Redaschi, and María Martín. FAIR adoption, assessment and challenges at UniProt. *Scientific Data*, 6, 12 2019.
- [5] Mark D. Wilkinson, Ruben Verborgh, Luiz Olavo Bonino da Silva Santos, Tim Clark, Morris A. Swertz, Fleur D.L. Kelpin, Alasdair J.G. Gray, Erik A. Schultes, Erik M. van Mulligen, Paolo Ciccicarese, Arnold Kuzniar, Anand Gavai, Mark Thompson, Rajaram Kaliyaperumal, Jerven T. Bolleman, and Michel Dumontier. Interoperability and FAIRness through a novel combination of Web technologies. *PeerJ Computer Science*, 3:e110, April 2017.

¹²<https://thehyve.nl/>.

¹³<http://www.fairdata.systems/>.

¹⁴<https://eccenca.com/>.