

# 1 **La unión hace la fuerza: modelos de distribución de especies**

## 2 **integrando diferentes fuentes de datos**

3 Javier Fernández-López<sup>1,2</sup>, Pelayo Acevedo<sup>3</sup>, Olivier Gimenez<sup>1</sup>

4 (1) CEFE, Université Montpellier, CNRS, EPHE, IRD, Montpellier (Francia)

5 (2) Universidad Complutense de Madrid, Madrid (España)

6 (3) IREC, CSIC-UCLM-JCCM, Ciudad Real (España)

7 Autor para correspondencia: Javier Fernández-López ([javfer05@ucm.es](mailto:javfer05@ucm.es))

## 8 **Palabras clave**

9 Proceso de Puntos de Poisson; modelos jerárquicos; verosimilitud conjunta; NIMBLE

## 10 **Keywords**

11 Poisson Point Process; hierarchical models; joint-likelihood; NIMBLE

12 La distribución y la abundancia de las especies son dos métricas básicas implicadas directamente en  
13 la gestión y la conservación de las poblaciones silvestres (Tellería, 2012). En las últimas décadas, los  
14 modelos de distribución (SDM por sus siglas en inglés) se han convertido en una de las herramientas  
15 más utilizadas para el estudio de la ocupación (presencia/ausencia) y abundancia (número de  
16 individuos) de las especies (Elith y Leathwick, 2009). Estos modelos relacionan la presencia o  
17 abundancia de individuos con una serie de covariables predictoras, permitiendo i) conocer la relación  
18 de esas covariables con las poblaciones y ii) realizar predicciones espacialmente explícitas sobre la  
19 ocupación o abundancia de las especies en función de esas covariables. Entre las causas que han  
20 motivado el auge de estas herramientas se podrían mencionar el aumento de la disponibilidad de  
21 cartografía ambiental (Hijmans et al., 2005), el desarrollo de nuevos algoritmos para ajustar este tipo  
22 de modelos (Phillips et al., 2006) o la aparición de numerosas plataformas con bases de datos  
23 accesibles sobre la presencia o abundancia de las especies (p.ej. [GBIF](https://www.gbif.org/)). Recientemente, el  
24 crecimiento y diversidad de estas plataformas ha motivado la investigación de las aproximaciones  
25 analíticas para la integración de diferentes fuentes de datos en un único modelo que refleje la  
26 distribución o abundancia de las especies (*integrated species distribution models*, en adelante iSDM).  
27 Estas metodologías pueden ser de gran utilidad en programas de monitoreo de poblaciones a gran  
28 escala o en especies para las que no se dispone de mucha información, pues permiten combinar  
29 unos pocos muestreos sistemáticos de gran calidad y llevados a cabo por especialistas con bases de

30 datos más oportunistas, como las procedentes de la ciencia ciudadana (Miller et al., 2019). El uso de  
31 la verosimilitud conjunta (en adelante *joint-likelihood*) es una de las metodologías más utilizadas para  
32 ajustar este tipo de modelos y se asienta principalmente en dos marcos teórico-prácticos: el Proceso  
33 de Puntos de Poisson y la modelización jerárquica.

## 34 **Proceso de Puntos de Poisson**

35 El Proceso de Puntos de Poisson es un modelo matemático que estudia la distribución de elementos  
36 (puntos) en un espacio continuo, asumiendo que siguen una distribución de Poisson cuyo parámetro  
37  $\lambda_i$  significaría el número esperado de puntos por unidad de área en la localidad  $i$ . Además de resultar  
38 una forma natural de representar individuos en el territorio (puntos en el espacio), este proceso tiene  
39 dos propiedades que lo hacen muy atractivo a la hora de utilizarlo como modelo subyacente en el  
40 estudio de la distribución y abundancia de las especies. La primera es que partiendo de una  
41 distribución de individuos (puntos en un espacio continuo), es muy fácil e intuitivo *discretizar* ese  
42 espacio (colocando encima una cuadrícula, por ejemplo) y obtener dos métricas de gran interés: la  
43 abundancia (conteo del número de puntos por cuadrícula) y la ocupación (cuadrículas con al menos  
44 un punto) de la especie (Figura 1). Esto significa que la ocupación y la abundancia pueden ser  
45 entendidas como realizaciones discretizadas del mismo proceso de puntos que ocurre en un espacio  
46 continuo y que, por tanto, están relacionadas entre sí. Esta propiedad permite combinar distintos  
47 tipos de datos (presencias/ausencias, solo presencias y conteos) en un mismo marco estadístico. La  
48 segunda propiedad tiene que ver con la resolución. El proceso de puntos ocurre, tal y como sucede  
49 con la distribución de individuos en la naturaleza, en un espacio continuo y por lo tanto es  
50 independiente de la resolución. Son las métricas *abundancia* y *ocupación* las que son dependientes  
51 de la resolución, ya que previamente hemos discretizado el espacio en cuadrículas cuyo tamaño va a  
52 condicionar el número o la presencia/ausencia de puntos en su interior (Figura 1). Por ello, utilizar el  
53 Proceso de Puntos de Poisson como modelo subyacente nos permite reconciliar diferentes bases de  
54 datos a diferentes resoluciones, ya que asumimos que todas provienen de la misma distribución de  
55 individuos en un espacio continuo.

## 56 **Modelos jerárquicos**

57 La modelización jerárquica es un marco estadístico que se caracteriza por diferenciar entre una  
58 variable latente (o variable de estado), que suele ser el objeto de nuestro interés pero que no puede  
59 ser estudiada directamente, y unas observaciones, que son muestreos de esa variable latente y que  
60 aunque dependen de ella, emergen a través de un proceso observacional y por lo tanto pueden estar  
61 distorsionadas (Kéry y Royle, 2015). En nuestro caso, la variable latente es la abundancia o la  
62 ocupación de la especie, mientras que las observaciones pueden ser muestreos de conteos directos,  
63 muestreos de detecciones/no-detecciones, registros de presencias oportunistas o cualquier otra  
64 fuente de datos (Fernández-López et al., 2022). Todas estas observaciones estarán relacionadas con

65 la variable latente, pero además se verán afectadas por un proceso observacional en el que influirán  
66 otros factores. Uno de los factores que se intenta controlar más a menudo es la detectabilidad  
67 imperfecta (falsos negativos: no detectar un individuo/especie cuando en realidad sí está presente en  
68 un determinado sitio), aunque en el marco de la modelización jerárquica se pueden acomodar  
69 igualmente otras situaciones, como el sesgo de muestreo o los falsos positivos (Louvrier et al., 2018).  
70 La ventaja de los modelos jerárquicos es que permiten incluir aquellas covariables que afectan a los  
71 dos procesos (ecológico para la variable latente y observacional para la fuente de datos) de forma  
72 independiente y mediante ajustes habituales de los modelos lineales generalizados (Figura 2).

### 73 **Modelos de distribución de especies integrados: *joint-likelihood***

74 La aproximación de *joint-likelihood* para el ajuste de iSDM utiliza como modelo subyacente el proceso  
75 de puntos de Poisson, siendo la variable latente el número de individuos por unidad de área (Miller  
76 et al., 2019). Partiendo de este modelo unificador, se superponen los diferentes juegos de datos con  
77 sus respectivos procesos (submodelos) observacionales, pero compartiendo una misma variable  
78 latente de la cual provienen (Figura 2). Al compartir la misma variable latente, las fuentes de datos  
79 comparten información sobre el proceso ecológico que las origina, teniendo en cuenta a su vez los  
80 diferentes procesos observacionales específicos de los cuales emanan cada una de ellas. Es  
81 importante destacar que para ajustar un iSDM mediante *joint-likelihood* se requiere una comprensión  
82 profunda de cada uno de los procesos observacionales que originan las diferentes fuentes de datos.  
83 Si alguno de los submodelos observacionales no es adecuado por cualquier motivo, la información  
84 compartida por esa fuente de datos será errónea, por lo que el resultado del modelo integrado estará  
85 así mismo sesgado (Ahmad Suhaimi et al., 2021).

86 La Figura 3 representa de forma esquemática los componentes y resultados de un ejemplo de iSDM  
87 a partir de datos simulados (ver Apéndice 1 para el código completo). En este modelo se combinaron  
88 dos fuentes de datos: i) una resultante de un muestreo en el que se visitaron 4 veces un total de 10  
89 localidades, contabilizando todos los individuos que fueron capaces de detectarse (“conteos  
90 repetidos”); y ii) una base de datos de ciencia ciudadana que contenía un total de 223 registros de  
91 presencia oportunistas (“registros oportunistas”). A partir de estos datos, se ajustaron tres modelos,  
92 uno con cada una de las fuentes de datos por separado y un tercer modelo integrando ambas  
93 fuentes. Para todos los modelos se siguió una estrategia de modelización jerárquica similar, en la que  
94 se asumió un Proceso de Puntos de Poisson como modelo subyacente para la variable latente  
95 (abundancia de individuos), y se ajustó un submodelo observacional para cada una de las fuentes de  
96 datos siguiendo el esquema de la Figura 2. Todos los modelos se ajustaron mediante inferencia  
97 Bayesiana en NIMBLE (Valpine et al., 2017). En la Figura 3 se muestran las fuentes de datos, las  
98 covariables predictoras y los principales resultados de los análisis realizados. Los patrones de  
99 abundancia obtenidos a partir de los diferentes modelos fueron similares a la abundancia simulada  
100 (Figura 3, segunda fila), mostrando una mayor linealidad para el modelo a partir de registros  
101 oportunistas y el iSDM. Sin embargo, el iSDM obtuvo estimas de abundancia totales más cercanas a

102 la abundancia simulada que el resto de modelos (N abundancia simulada = 25,801; N Conteos SDM  
103 = 14,656.; N Registros oportunistas = 17,107; N iSDM = 18,003). La precisión de los parámetros  
104 estimados fue mayor en el caso del iSDM (Apéndice 1).

105 Los modelos integrados permiten combinar la información de calidad obtenida a partir de muestreos  
106 sistemáticos, pero que normalmente son escasos y de ámbito local por su elevado coste, con otro  
107 tipo de datos oportunistas, como los derivados de la ciencia ciudadana, que suelen ser muy  
108 abundantes y ampliamente distribuidos pero que a menudo se ven gravemente afectados por sesgos  
109 de diferente tipo (Dorazio, 2014; Isaac et al., 2020; Merow et al., 2022). El desarrollo de los modelos  
110 de distribución integrados permitirá el aprovechamiento de las diferentes bases de datos de una  
111 forma más eficiente, cobrando especial importancia en programas de monitorización a gran escala y  
112 en aquellas especies menos estudiadas o con mayores problemas de conservación (Tellería, 2012).

### 113 **Contribución de los autores**

114 Javier Fernández-López: Conceptualización, análisis formal y redacción - borrador original. Pelayo  
115 Acevedo: Conceptualización y redacción - revisión y edición. Olivier Gimenez: Conceptualización,  
116 validación y redacción - revisión y edición.

### 117 **Agradecimientos**

118 Queremos agradecer a Verónica Cruz, Julen Astigarraga, Elena Quintero y en general al Grupo de  
119 Ecoinformática de la Asociación Española de Ecología Terrestre por organizar y promover los  
120 “Seminarios Ecoinformáticos”, a partir de uno de los cuales se originó esta nota. Las conversaciones  
121 con Valentin Lauret motivaron reflexiones que ayudaron a una mejor evaluación y comunicación de  
122 los modelos de esta nota. JF-L está actualmente financiado por un contrato postdoctoral Margarita  
123 Salas (fondos de la Unión Europea – NextGenerationEU) a través de la Universidad Complutense de  
124 Madrid.

125 REFERENCIAS

- 126 Ahmad Suhaimi, S.S., Blair, G.S., Jarvis, S.G. 2021. Integrated species distribution models: A  
127 comparison of approaches under different data quality scenarios. *Diversity and Distributions* 27:  
128 1066-1075.
- 129 Dorazio, R.M. 2014. Accounting for imperfect detection and survey bias in statistical analysis of  
130 presence-only data. *Global Ecology and Biogeography* 23: 1472-1484.
- 131 Elith, J., Leathwick, J.R. 2009. Species distribution models: ecological explanation and prediction  
132 across space and time. *Annual Review of Ecology, Evolution and Systematics* 40: 677-697.
- 133 Fernández-López, J., Blanco-Aguilar, J.A., Vicente, J., Acevedo, P. 2022. Can we model distribution of  
134 population abundance from wildlife–vehicles collision data? *Ecography* 2022: e06113.
- 135 Hijmans, R.J., Cameron, S.E., Parra, J.L., Jones, P.G., Jarvis, A. 2005. Very high resolution  
136 interpolated climate surfaces for global land areas. *International Journal of Climatology: A  
137 Journal of the Royal Meteorological Society* 25: 1965-1978.
- 138 Isaac, N.J., Jarzyna, M.A., Keil, P., Dambly, L.I., Boersch-Supan, P.H., Browning, E., Freeman, S.N.  
139 et al. 2020. Data integration for large-scale models of species distributions. *Trends in ecology &  
140 evolution* 35: 56-67.
- 141 Kéry, M., Royle, J.A. 2015. *Applied Hierarchical Modeling in Ecology: Analysis of Distribution,  
142 Abundance and Species Richness in R and BUGS*. Elsevier, Londres, Reino Unido.
- 143 Louvrier, J., Chambert, T., Marboutin, E., Gimenez, O. 2018. Accounting for misidentification and  
144 heterogeneity in occupancy studies using hidden Markov models. *Ecological modelling* 387: 61-  
145 69.
- 146 Merow, C., Galante, P.J., Kass, J.M., Aiello-Lammens, M.E., Babich Morrow, C., Gerstner, B.E.,  
147 Grisales Betancur, V. et al. 2022. Operationalizing expert knowledge in species' range  
148 estimates using diverse data types. *Frontiers of Biogeography*.
- 149 Miller, D.A., Pacifici, K., Sanderlin, J.S., Reich, B.J. 2019. The recent past and promising future for  
150 data integration methods to estimate species' distributions. *Methods in Ecology and Evolution*  
151 10: 22-37.
- 152 Phillips, S.J., Anderson, R.P., Schapire, R.E. 2006. Maximum entropy modeling of species geographic  
153 distributions. *Ecological modelling* 190: 231-259.
- 154 Tellería, J.L. 2012. *Introducción a la conservación de las especies*. Hernández, V. (ed.), Tundra  
155 Ediciones, Valencia, España.

156 Valpine, P. de, Turek, D., Paciorek, C.J., Anderson-Bergman, C., Lang, D.T., Bodik, R. 2017.  
157 Programming with models: writing statistical algorithms for general model structures with  
158 NIMBLE. *Journal of Computational and Graphical Statistics* 26: 403-413.

159 PIES DE FIGURA

160 **Figura 1.** El mismo proceso de puntos (A y D) y sus métricas derivadas a diferentes resoluciones: B y  
161 C a 1x1 km; E y F a 4x4 km. Los gráficos B y E muestran el número de puntos por cuadrícula  
162 (abundancia). Los gráficos C y F muestran la presencia de al menos un punto en cada una de las  
163 cuadrículas (ocupación). Mientras que los patrones observados a partir de las métricas derivadas  
164 (abundancia B y E; ocupación C y F) varían dependiendo de la resolución de la malla, el patrón de  
165 puntos permanece constante (A y D).

166 **Figure 1.** The same point process (A and D) and their derived metrics at different resolutions: B and  
167 C at 1x1 km; E and F at 4x4 km. Graphs B and E show the number of points per grid (abundance).  
168 Graphs C and F show the presence of at least one point in each of the grids (occupancy). While the  
169 patterns observed from the derived metrics (abundance B and E; occupancy C and F) vary depending  
170 on the resolution of the grid, the pattern of points remains constant (A and D).

171 **Figura 2.** Esquema de un modelo jerárquico integrando dos fuentes de datos diferentes (ver  
172 Apéndice 1). A) La variable latente o de estado (abundancia de individuos) puede ser modelizada  
173 como un Proceso de Puntos de Poisson discretizado, cuya función de intensidad ( $\lambda$ ) vendrá  
174 determinada por dos covariables (cobertura de bosque y altitud). B) La primera fuente de datos,  
175 conteos sistemáticos en 15 localidades repetidos cuatro veces por cada localidad, puede ser  
176 modelizada con una distribución Binomial (detección - no detección de cada uno de los individuos en  
177 cada conteo y localidad). C) La segunda fuente de datos, registros oportunistas a partir de ciencia  
178 ciudadana, puede modelizarse como un Proceso de Puntos de Poisson adelgazado con una  
179 distribución de Bernoulli (para cada celda, detección - no detección de al menos un individuo).

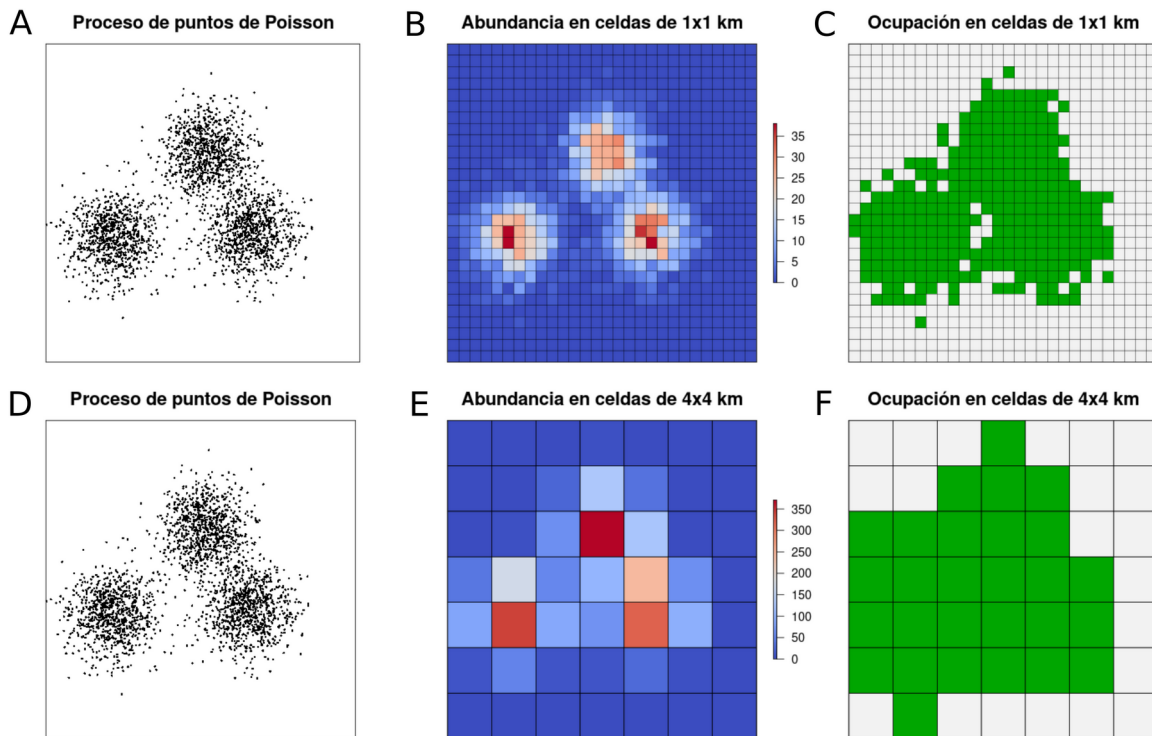
180 **Figure 2.** Diagram of a hierarchical model integrating two different data sources (see Appendix 1). A)  
181 The latent or state variable (abundance of individuals) can be modeled as a discretized Poisson Point  
182 Process, whose intensity function ( $\lambda$ ) will be determined by two covariates (forest cover and altitude).  
183 B) The first source of data, systematic counts in 15 locations repeated four times at each location, can  
184 be modeled with a Binomial distribution (detection - non-detection of each individual in each count and  
185 location). C) The second data source, opportunistic records from citizen science, can be modeled as a  
186 thinned Poisson Point Process with a Bernoulli distribution (for each cell, detection - non-detection of  
187 at least one individual).

188 **Figura 3.** Fuentes de datos, covariables predictoras y resultados del ejemplo desarrollado en el  
189 Apéndice 1. Se muestra el patrón de distribución de la abundancia simulada (Abundancia), y las  
190 predicciones para cada uno de los modelos: a partir de conteos repetidos (Conteos SDM), a partir de  
191 registros oportunistas (Registros SDM) y el modelo integrando ambas bases de datos (iSDM).

192 **Figure 3.** Data sources, predictor covariates and results of the example developed in Appendix 1. We  
193 show the simulated abundance distribution pattern (Abundancia), and the predictions for each of the

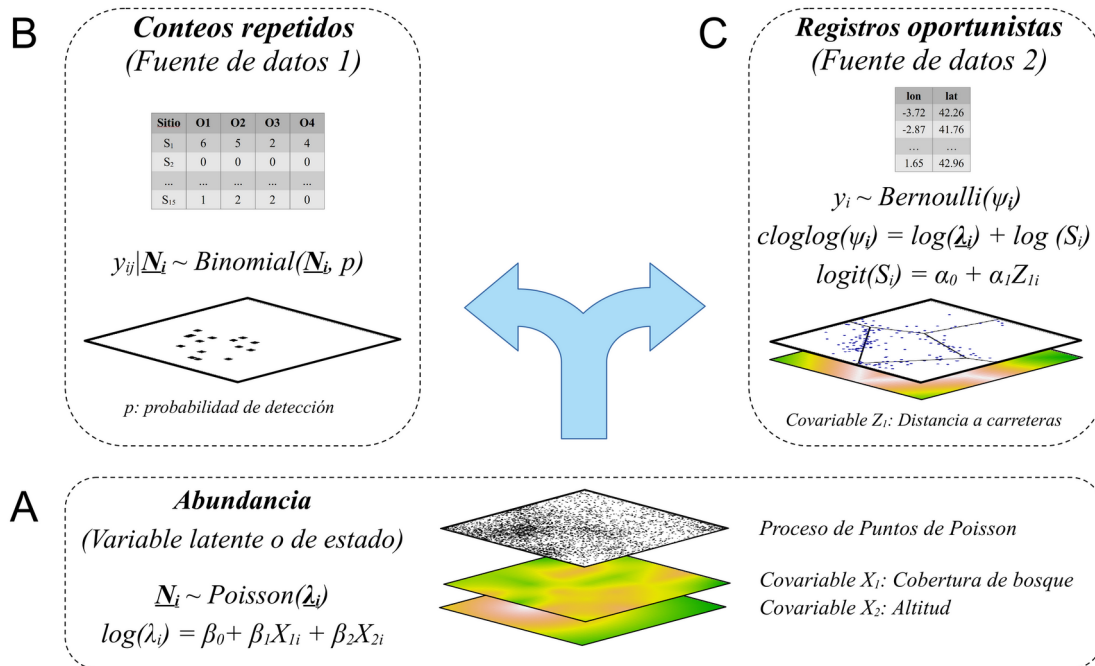
194 models: from repeated counts (Conteos SDM), from opportunistic records (Registros SDM) and the  
195 model integrating both databases (iSDM).





197

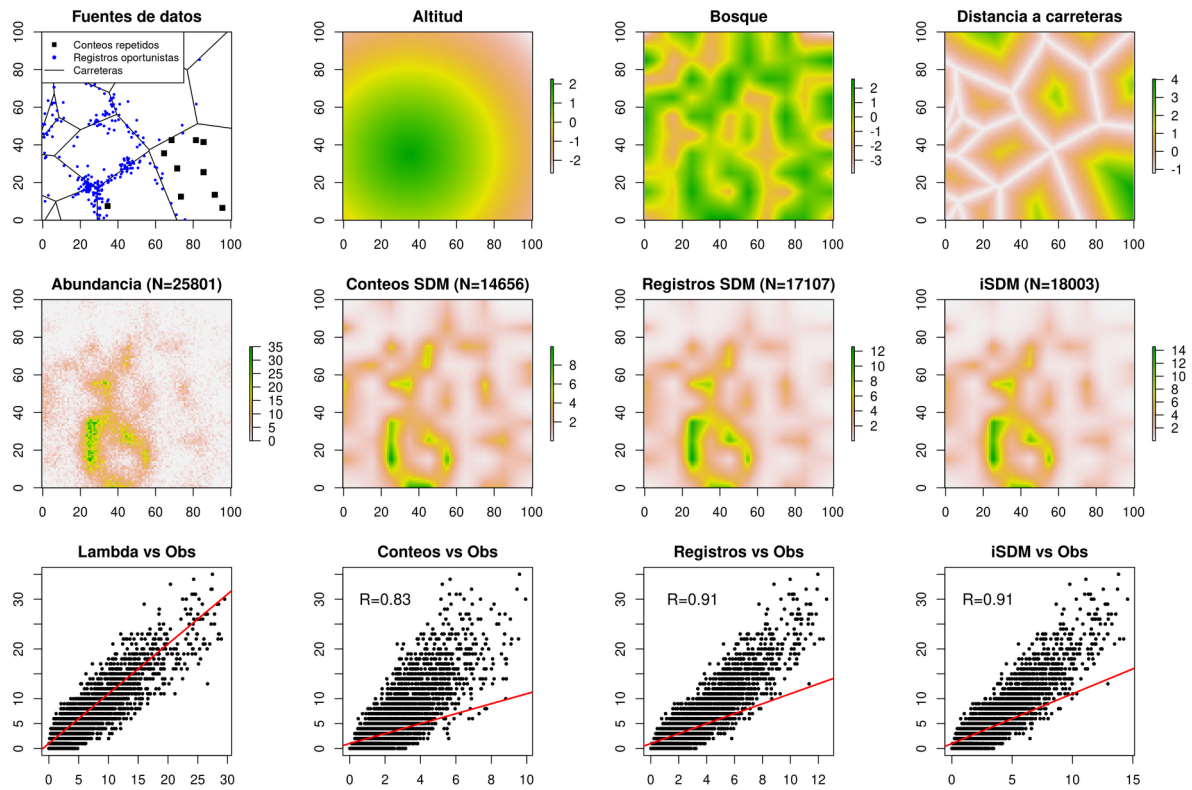
198 *Figura 1. El mismo proceso de puntos (A y D) y sus métricas derivadas a diferentes resoluciones: B y*  
 199 *C a 1x1 km; E y F a 4x4 km. Los gráficos B y E muestran el número de puntos por cuadrícula*  
 200 *(abundancia). Los gráficos C y F muestran la presencia de al menos un punto en cada una de las*  
 201 *cuadrículas (ocupación). Mientras que los patrones observados a partir de las métricas derivadas*  
 202 *(abundancia B y E; ocupación C y F) varían dependiendo de la resolución de la malla, el patrón de*  
 203 *puntos permanece constante (A y D).*



205

206 *Figura 2. Esquema de un proceso de modelado jerárquico integrando dos fuentes de datos diferentes*  
 207 *(ver Apéndice 1). A) La variable latente o de estado (abundancia de individuos) puede ser*  
 208 *modelizada como un Proceso de Puntos de Poisson discreto, cuya función de intensidad (lambda)*  
 209 *vendrá determinada por una serie de covariables (cobertura de bosque y altitud). B) La primera*  
 210 *fuentes de datos, conteos sistemáticos de 15 localidades repetidos cuatro veces en cada localidad,*  
 211 *puede ser modelizada con una distribución Binomial (detección - no detección de cada uno de los*  
 212 *individuos en cada conteo y localidad). C) La segunda fuente de datos, registros oportunistas a partir*  
 213 *de ciencia ciudadana, puede modelizarse como un Proceso de Puntos de Poisson adelgazado con*  
 214 *una distribución de Bernoulli (para cada celda, detecto - no detecto al menos un individuo).*

215 FIGURA 3



216

217 *Figura 3. Fuentes de datos, covariables predictoras y resultados del ejemplo desarrollado en el*  
 218 *Apéndice 1. Se muestra el patrón de distribución de la abundancia simulada (Abundancia), y las*  
 219 *predicciones para cada uno de los modelos: a partir de conteos repetidos (Conteos SDM), a partir de*  
 220 *registros oportunistas (Registros SDM) y el modelo integrando ambas bases de datos (iSDM).*