

Tres sistemas de policía predictiva en España: VioGén, RisCanvi y VeriPol

Evaluación desde una perspectiva de derechos humanos

Lucía Martínez Garay (coord.)



Tres sistemas de policía predictiva en España: VioGén, RisCanvi y VeriPol

Evaluación desde una perspectiva
de derechos humanos

Tres sistemas de policía predictiva en España: VioGén, RisCanvi y VeriPol

Evaluación desde una perspectiva
de derechos humanos

Lucía Martínez Garay, coord.

Colección: Desarrollo Territorial
Papers, 10

Dirección de la colección: Maria Dolores Pitarch

Consejo de dirección: Josep Vicent Boira, Sacramento Pinazo, Joan Romero i Ana Sales



Esta obra está bajo una Licencia Creative Commons Reconocimiento-NoComercial-SinObraDerivada 4.0 Internacional.

Esta publicación no puede ser reproducida, ni total ni parcialmente, ni registrada en, o transmitida por, un sistema de recuperación de información, de ninguna forma ni por ningún medio, sea fotomecánico, fotoquímico, electrónico, por fotocopia o por cualquier otro, sin el permiso de la editorial. Diríjase a CEDRO (Centro Español de Derechos Reprográficos, www.cedro.org) si necesita fotocopiar o escanear algún fragmento de esta obra.

Título original:

Three predictive policing approaches in Spain: VioGén, RisCanvi and VeriPol. Assessment from a human rights perspective

© Universitat de València, 2024

© Del texto: los autores y las autoras, 2025

© De esta edición: Universitat de València, 2025

Publicacions de la Universitat de València

puv.uv.es

publicacions@uv.es

Coordinación editorial: Amparo Jesús-Maria Romero

Corrección y maquetación: Letras y Píxeles, S. L.

Traducción: traductor automático, supervisado por los autores y las autoras

Diseño de la cubierta: Celso Hernández de la Figuera

ISBN papel: 978-84-1118-659-9

ISBN OA: 978-84-1118-660-5

DOI: <https://doi.org/10.7203/PUV-OA-9788411186605>

Edición digital

ÍNDICE

| | |
|---|-----------|
| CONTRIBUCIONES | 11 |
| AGRADECIMIENTOS | 12 |
| LISTADO DE ABREVIATURAS | 13 |
| RESUMEN EJECUTIVO | 15 |
| 1. Génesis del informe y elección de las herramientas estudiadas | 16 |
| 2. Condicionantes y límites de la realización del informe | 16 |
| 3. Conclusiones generales aplicables a las tres herramientas analizadas | 18 |
| 4. Principales conclusiones respecto de VioGén | 21 |
| 5. Principales conclusiones respecto de RisCanvi..... | 22 |
| 6. Principales conclusiones respecto de VeriPol | 24 |
| VIOGÉN | 25 |
| 1. Descripción general | 25 |
| 2. Datos | 27 |
| 3. Eficacia del enfoque policial predictivo para prevenir la delincuencia | 51 |
| 4. Evaluación de impacto sobre los derechos humanos..... | 66 |
| 5. Nuevos desarrollos sobre VeriPol entre enero y abril de 2023..... | 73 |
| 6. Fuentes..... | 76 |
| RISCANVI..... | 79 |
| 1. Descripción general | 79 |
| 2. Datos | 85 |
| 3. Eficacia del enfoque policial predictivo para prevenir la delincuencia | 113 |
| 4. Evaluación del impacto sobre los derechos humanos | 118 |
| 5. Fuentes..... | 133 |

| | |
|---|------------|
| VERIPOL | 135 |
| 1. Descripción general | 135 |
| 2. Datos | 137 |
| 3. Eficacia del enfoque policial predictivo para prevenir la delincuencia | 145 |
| 4. Evaluación del impacto sobre los derechos humanos | 161 |
| 5. Fuentes..... | 165 |
| ANEXOS SOBRE VIOGÉN..... | 167 |
| Anexo I: dificultades para obtener información y datos sobre VeriPol..... | 167 |
| Anexo II: Los factores de riesgo de VioGén | 172 |
| Anexo III: Medidas de protección policial para cada nivel de riesgo según la Instrucción 4/2019 de la Secretaría de Estado de Seguridad..... | 179 |
| ANEXOS SOBRE RISCANVI | 187 |
| 1. Riscanvi completo..... | 187 |
| 2. <i>Riscanvi Screening</i> | 209 |
| ADENDA: PREGUNTAS Y RESPUESTAS..... | 213 |
| ADENDA AL INFORME «TRES ENFOQUES POLICIALES PREDICTIVOS EN ESPAÑA»:..... | 213 |
| 1. Preguntas específicas sobre VioGén | 213 |
| 2. Preguntas específicas sobre RisCanvi..... | 215 |
| 3. Preguntas específicas sobre VeriPol..... | 228 |

Informe realizado por:

Lucía Martínez Garay. Investigadora responsable. Departament de Dret Penal e Institut Universitari d'Investigació en Criminologia i Ciències Penals, Universitat de València.

Andrés Boix Palop. Departament de Dret Administratiu i Processal, Universitat de València.

Álvaro Briz Redón. Departament d'Estadística i Investigació Operativa, Universitat de València.

Fernando Flores Giménez. Departament de Dret Constitucional, Ciència Política i de l'Administració e Institut de Drets Humans, Universitat de València.

Andrea García Ortiz. Institut Universitari d'Investigació en Criminologia i Ciències Penals, Universitat de València.

Mireia Molina Sánchez. Departament de Dret Administratiu i Processal, Universitat de València.

Francisco Montes Suay. Departament d'Estadística i Investigació Operativa, Universitat de València.

Adrián Palma Ortigosa. Departament de Dret Administratiu i Processal, Universitat de València.

Alfred Peris Manguillot. Institut de Matemàtica Pura i Aplicada, Universitat Politècnica de València.

Alba Soriano Arnanz. Departament de Dret Administratiu i Processal, Universitat de València.

Financiación y conflictos de intereses

El origen de este trabajo se encuentra en un encargo de una organización no gubernamental de derechos humanos para estudiar las herramientas policiales predictivas existentes en España, organización que también proporcionó apoyo financiero. Este informe es fruto de la actividad investigadora desarrollada en el marco del proyecto de investigación «Algorithmic Law» (CIPROM/2024/16) de la Universitat de València, del que forman parte los siguientes autores: Andrés Boix Palop, Lucía Martínez Garay, Alba Soriano Arnanz, Adrián Palma Ortigosa y Mireia Molina Sánchez; en el proyecto de investigación PID2021-123441NB-I00 (financiado por MCIN/AEI/ 10.13039/501100011033 y por «FEDER A way of making Europe», del que son miembros los siguientes autores: Lucía Martínez Garay y Andrea García Ortiz; y por el proyecto de investigación PID2023-152781NB-I00 «Administrative Planning» (PlanAdm), del que forman parte los siguientes autores: Andrés Boix Palop, Alba Soriano Arnanz y Adrián Palma Ortigosa.

Las opiniones expresadas en este informe son exclusivamente las de los autores y no representan las de las instituciones a las que pertenecen, ni las de las instituciones que encargaron o financiaron la investigación.



Limitaciones

Este informe se cerró en noviembre de 2022, pero diversas razones técnicas y personales han retrasado su publicación.

Con posterioridad a esa fecha solo se han incluido de forma muy fragmentaria algunas informaciones y noticias aparecidas sobre VioGén y RisCanvi hasta abril de 2023, así como algunas referencias muy breves al informe de auditoría publicado sobre RisCanvi en 2024.

La compatibilidad de las tres herramientas de policía predictiva con los derechos humanos se ha evaluado en relación con la normativa vigente en España y Europa a la fecha de cierre del informe. Por ello, en este informe no se evalúa la conformidad de las herramientas con las obligaciones establecidas en el Reglamento (UE) 2024/1689 del Parlamento Europeo y del Consejo, de 13 de junio de 2024, por el que se establecen normas armonizadas en materia de inteligencia artificial en la Unión Europea.

Contribuciones

Lucía Martínez Garay ha sido la investigadora responsable y ha dirigido y coordinado el informe. Participó en el análisis de datos cuantitativos de VioGén y RisCanvi, entrevistó a agentes de policía, abogados, médicos forenses y a un magistrado, fue responsable de la redacción preliminar del capítulo de VioGén, participó en la redacción preliminar de los capítulos de RisCanvi y VeriPol, y ha supervisado toda la redacción final del informe completo.

Francisco Montes Suay y Álvaro Briz Redón realizaron el análisis cuantitativo de los datos sobre RisCanvi y colaboraron en el análisis de algunas de las muestras de VioGén.

Alfred Peris Manguillot realizó el análisis cuantitativo de los datos públicamente disponibles sobre VioGén y ayudó a revisar el análisis de las muestras de VioGén.

Alba Soriano Arnanz y Adrián Palma Ortigosa analizaron las sentencias de las audiencias provinciales sobre RisCanvi. Se encargaron de la redacción preliminar de los capítulos de RisCanvi y VeriPol, y entrevistaron a los abogados. Adrián Palma Ortigosa entrevistó a algunos agentes de policía sobre VeriPol y Alba Soriano Arnanz ha traducido la mayor parte del informe al inglés y ha revisado la traducción realizada por los demás miembros del equipo.

Andrea García Ortiz y Mireia Molina Sánchez participaron en el análisis de datos cuantitativos sobre VioGén y entrevistaron a los jueces. Mireia Molina Sánchez analizó las sentencias de las audiencias provinciales sobre VioGén, y Andrea García Ortiz entrevistó a una trabajadora de la Oficina de Asistencia a la Víctima del Delito.

Andrés Boix Palop coordinó las entrevistas a los jueces, realizó algunas de ellas, supervisó las partes relativas a las obligaciones de rendición de cuentas de cada herramienta y revisó la redacción final de todo el informe.

Fernando Flores Jiménez facilitó el contacto con la Secretaría de Estado de Seguridad y revisó la redacción final de todo el informe.

Todas las conclusiones han sido debatidas y asumidas por todos los autores. Las opiniones expresadas en este informe son las de sus autores y no representan las de las instituciones a las que pertenecen, ni tampoco las de las instituciones que encargaron o financiaron la investigación.

Agradecimientos

Queremos agradecer a las siguientes personas su ayuda o contribución a este informe:

Carlos Castillo y Marzieh Karimi Haghghi, de la Universitat Pompeu Fabra, por los cálculos que hicieron para nosotros con sus muestras de datos de RisCanvi.

Virginia Álvarez, por su apoyo durante la elaboración del informe.

Verónica Gisbert Gracia, de la Universitat de València, por sus orientaciones y consejos sobre metodologías cualitativas y en la preparación de las entrevistas.

Antonio Andrés Pueyo, de la Universitat de Barcelona, por facilitar información sobre la descripción detallada de algunos de los factores de riesgo de RisCanvi.

El Tribunal Superior de Justicia de la Comunidad Valenciana, la Asociación de Psicología Forense de la Administración de Justicia, la Secretaría de Estado de Seguridad del Ministerio del Interior, la Unidad de Planificación y Coordinación Estratégica de la Policía Nacional, las jefas de las Unidades de Valoración Forense de los Institutos de Medicina Legal y Ciencias Forenses de Valencia y Alicante, la persona del Departament de Justícia de la Generalitat de Catalunya que contestó a nuestros correos electrónicos, los jueces, policías, abogados, personal de la Oficina de Asistencia a las Víctimas del Delito de Valencia, y todas las personas que estuvieron dispuestas a darnos información sobre alguna de las herramientas de evaluación del riesgo o nos facilitaron el contacto con personas que podían hacerlo.

Listado de abreviaturas

| | |
|--------------|--|
| AEPD | Agencia Española de Protección de Datos |
| CGPJ | Consejo General del Poder Judicial |
| RGPD | Reglamento General de Protección de Datos |
| HCR-20 | <i>Historical Clinical Risk Management–20</i> |
| IMLCF | Instituto de Medicina Legal y Ciencias Forenses |
| JVM | Jueces de Violencia sobre la Mujer |
| LSI-R™ | <i>Level of Service Inventory–Revised™</i> |
| TM | Trastorno mental |
| PCL-R | Lista de Comprobación de Psicopatía - Revisada |
| JVP | Juez de vigilancia penitenciaria |
| REVI | Escala de riesgo de reincidencia violenta de RisCanvi |
| RisCanvi-C | RisCanvi completo |
| RisCanvi-S | RisCanvi <i>Screening</i> |
| SARA | <i>Spousal Assault Risk Assessment Guide</i> |
| SIDENPOL | Sistema de Denuncias Policiales de la Policía Nacional |
| SIPC | Sistema de Información Penitenciaria de Cataluña |
| SOS-RisCanvi | Equipo de apoyo, orientación y seguimiento en la aplicación del RisCanvi |
| SES | Secretaría de Estado de Seguridad |
| SVR-20 | <i>Sexual Violence Risk–20</i> |
| VFR | Valoración forense del riesgo |
| VPER | Valoración policial de la evolución del riesgo |
| VPR | Valoración policial del riesgo |

RESUMEN EJECUTIVO

El informe que se enlaza a continuación, *Three predictive policing approaches in Spain: Viogén, RisCanvi and VeriPol (Assessment from a human rights perspective)*, ha sido elaborado por un grupo de profesores de la Universitat de València y la Universitat Politècnica de València, especialistas en distintas ramas del derecho (penal, administrativo y constitucional), la estadística y la matemática, bajo la coordinación de la Dra. Lucía Martínez Garay, profesora de Derecho Penal de la Universitat de València. El trabajo evalúa tres herramientas algorítmicas (con uso de inteligencia artificial, IA, en una de ellas) empleadas por la Administración española para realizar análisis predictivos en los ámbitos policiales y penitenciarios, desde la perspectiva de su respeto a los derechos fundamentales de los ciudadanos que son objeto de las evaluaciones para las que estas herramientas son empleadas. Los objetivos de este informe han sido: 1) analizar la capacidad predictiva de los tres sistemas y su eficacia para alcanzar el objetivo que pretenden; 2) analizar su grado de transparencia, tanto para la ciudadanía en general como para las personas directamente afectadas por su utilización; 3) identificar posibles sesgos estadísticos; 4) analizar si respetan el vigente marco jurídico europeo y español, tanto en materia de protección de datos como respecto del empleo de este tipo de modelos; y, finalmente, 5) valorar si los costes para los derechos humanos exceden lo que sería legítimo, necesario o proporcionado en relación con el beneficio actual o probable en términos de prevención del delito. Aunque estos han sido los objetivos comunes para las tres herramientas, el análisis se ha tenido que adaptar a las características de cada una de ellas y, sobre todo, al tipo y a la cantidad de datos e información de los que hemos podido disponer, que ha sido dispar en cada una de las tres. Asimismo, en el análisis se han combinado métodos cuantitativos y cualitativos, de nuevo de forma distinta para cada herramienta según la información a la que nos ha sido posible acceder.

1. GÉNESIS DEL INFORME Y ELECCIÓN DE LAS HERRAMIENTAS ESTUDIADAS

El origen de este informe es un encargo de una ONG de ámbito europeo e internacional que pretende realizar un estudio comparativo respecto de los riesgos que las herramientas de policía predictiva utilizadas en distintos estados europeos pueden plantear para los derechos fundamentales de los ciudadanos. El presente trabajo se refiere a la parte del análisis sobre las herramientas de uso predictivo empleadas en España, y ha sido realizado, en el marco de este encargo, a partir de un contrato de transferencia del conocimiento con la Universitat de València, bajo la dirección y coordinación de Lucía Martínez Garay. Sus resultados son, además, y en parte, fruto de la tarea de investigación realizada por la mencionada profesora y el resto del equipo en varios proyectos de investigación («Algorithmical Law», PROMETEU/2021/009; PID2021-123441NB-I00, financiado por MCIN/AEI/10.13039/501100011033/ y «ERDF A way of making Europe»; y «Digital economy regulation: equality guarantees provided by public powers and algorithmical tools», PID2019-108745GB-I00). En todo caso, los análisis, las conclusiones y las opiniones expresadas en este informe son únicamente los de sus autores y no representan los de las instituciones a las que pertenecen ni los de las que han financiado estas investigaciones, así como tampoco a la ONG que lo ha encargado.

Este informe no analiza todas las herramientas de esta índole que en estos momentos emplean las Administraciones españolas, sino únicamente tres de ellas, que seleccionó la ONG que encargó el informe a partir de un listado más amplio proporcionado por el equipo investigador. Fue también la ONG quien definió los aspectos que debían ser analizados en cada una. En concreto, se han estudiado dos herramientas algorítmicas utilizadas por la policía (VeriPol, empleada para la detección de posibles denuncias falsas de robo; y VioGén, utilizada para estimar los niveles de riesgo para las víctimas de violencia de género, a fin de decidir mejor las medidas de protección policial necesarias en cada caso) y una tercera que utiliza la Administración autonómica catalana para la gestión de sus prisiones (RisCanvi, que estima tanto el riesgo de reincidencia como el de otras conductas peligrosas por parte de la población reclusa).

2. CONDICIONANTES Y LÍMITES DE LA REALIZACIÓN DEL INFORME

Este informe fue finalizado y entregado en noviembre de 2022, aunque su publicación por parte de la ONG que lo encargó, junto al resto de informes de otros países europeos, aún no ha tenido lugar. No obstante, en esta versión se han añadido, aunque de manera sucinta y fragmentaria, ciertos datos e información sobre VioGén y RisCanvi obtenidos hasta abril de 2023, así como algunas referencias a la auditoría que las propias autoridades autonómicas catalanas encargaron respecto de RisCanvi y que ha sido publicada en el primer trimestre de 2024. Estas fechas han condicionado también el marco normativo de referencia respecto al cual se ha evaluado la compatibilidad de las herramientas con los derechos fundamentales: las regulaciones española y europea vigentes en materia de empleo de herramientas

automatizadas por los poderes públicos, legislación sobre modelos predictivos y normativa sobre protección de datos de carácter personal, vigente en 2022. Por esta razón, el informe no evalúa el grado de cumplimiento de las herramientas con las obligaciones que el Reglamento Europeo en materia de Inteligencia Artificial (AI Act), aprobado por el Parlamento Europeo el 13 de marzo de 2024 y publicado por el *Diario Oficial de la Unión Europea* en junio de 2024, establece respecto del uso de sistemas algorítmicos por parte de los poderes públicos. Téngase en cuenta, a estos efectos, que la limitada definición de «Inteligencia Artificial» que el artículo 3 de este reglamento contiene en su versión final, según la cual esta requiere de un sistema diseñado para funcionar no solo con diversos niveles de autonomía, sino que, además, sea capaz de mostrar capacidad de adaptación tras su despliegue, es decir, una capacidad de evolución/aprendizaje mínimamente autónomo, hace que de las tres herramientas analizadas solo VeriPol pueda entenderse que entra en la actualidad en sentido estricto en esta definición y debería, por tanto, cumplir con todas las previsiones del reglamento cuando entre en vigor para los usos de IA por parte de los poderes públicos (lo que no está previsto hasta dentro de unos años por el propio reglamento, en todo caso, lo que hace que este análisis sea más importante en cuanto a pautas y directrices que aplicar que de cumplimiento normativo estricto). Por lo demás, esto podría cambiar en un futuro, dado que hay, en el caso de VioGén, evoluciones en proceso dirigidas a incorporar moderadamente sistemas de IA para la mejora evolutiva de la herramienta a fin de incrementar su capacidad predictiva, y de que respecto de RisCanvi se han propuesto también modificaciones en este sentido (como resultado de la auditoría externa encargada por la Generalitat). En todo caso, la mayor parte de los juicios de valor sobre adaptación de las herramientas analizadas al marco jurídico y al cumplimiento de ciertas exigencias de transparencia y de auditabilidad se corresponden con el grueso de las exigencias que esta norma jurídica de derecho europeo ha incorporado para los supuestos de utilización de IA por parte de los poderes públicos. Pues entendemos que, si bien no como marco normativo cuyo cumplimiento sea ya directamente exigible a los poderes públicos, sí constituyen guías y directrices que, como mínimo, han de ser tenidas en consideración cuando emplean estas herramientas.

Por último, queremos señalar que la información sobre la que se ha elaborado el informe es en parte de acceso público, y en parte nos ha sido proporcionada o bien por las instituciones que gestionan cada una de las herramientas analizadas, o bien por personas que las utilizan o conocen por su actividad profesional. Respecto de esta segunda clase de información, solo hemos reproducido aquí los datos o extremos que no estaban expresamente cubiertos por un acuerdo de confidencialidad. También hemos mantenido en todo caso el anonimato de todas las fuentes empleadas cuando así lo han requerido y cuando su identificación personal no era relevante a efectos del informe realizado.

3. CONCLUSIONES GENERALES APLICABLES A LAS TRES HERRAMIENTAS ANALIZADAS

Las principales conclusiones que se extraen del informe, y que son predicables de las tres herramientas analizadas, son las siguientes.

- En general, todas ellas cumplen con el marco normativo vigente en España en materia de empleo de herramientas algorítmicas y sistemas automatizados para el apoyo a la acción administrativa, lo que en gran parte tiene que ver no con un diseño de estas particularmente riguroso y exigente en cuanto a sus garantías, sino con la muy limitada ambición del marco jurídico español vigente en esta materia. Igualmente, no se han detectado incumplimientos relevantes respecto de las obligaciones de derecho europeo y nacional en materia de protección de datos de carácter personal, lo que también es consecuencia de que resulta relativamente fácil cumplirlas en estos casos porque el marco no es particularmente exigente con los poderes públicos (que han de cumplir con las exigencias técnicas, pero ni siquiera siempre con la necesidad de consentimiento, que puede ser sustituido por una habilitación legal). En todo caso, en el informe se proponen mejoras en todos los casos en esta materia, que tienen que ver con una mejor auditoría y evaluación de los riesgos, que no parece haber sido realizada con la intensidad y profundidad exigida por el marco jurídico vigente, tanto en uno como en otro ámbito, cuando estamos ante sistemas que pueden tener una afección notable a los derechos de los ciudadanos. Asimismo, entendemos que hay que garantizar mejor el derecho de acceso de los ciudadanos a sus datos personales empleados por el sistema.
- El principal punto negro que comparten las tres herramientas, si bien en diversos grados, es la falta de transparencia, información y suministro de datos completos sobre su funcionamiento, que facilite una auditoría y control externos capaces de detectar eventuales sesgos discriminatorios o riesgos para los derechos de los ciudadanos, por una parte; y que permita, por otra, evaluar mejor la capacidad predictiva y su eficacia en la prevención del delito, sin tener que confiar para ello únicamente en los informes y evaluaciones de las propias administraciones públicas que las emplean o han diseñado. En relación con este punto, es importante, a nuestro juicio, diferenciar dos cuestiones. Una cosa es que el diseño, la implementación o incluso el control del funcionamiento de estos sistemas hayan sido rigurosos, realizados por especialistas, muchas veces en saludable colaboración entre las administraciones y las universidades –lo que es sin duda una primera garantía de que las herramientas están construidas con rigor y calidad desde el punto de vista técnico–, y otra bien distinta que a pesar de ello sigan siendo necesarios tanto la transparencia de cara a la ciudadanía como el control externo por parte de terceros no implicados en el diseño ni el funcionamiento diario de los sistemas. Lo primero, porque tanto el diseño de estas herramientas como el establecimiento de las pautas sobre cómo usarlas implican decisiones tanto técnicas como sobre qué fines deben ser priorizados, y qué costes se considera aceptable asumir para alcanzar aquellos fines. Decisiones que son por lo tanto, en última instancia,

de naturaleza política, y sobre las cuales los ciudadanos tienen derecho a estar informados, por razones de confianza de la ciudadanía en las herramientas y en los poderes públicos que las emplean, y para poder controlar la acción de estos. Especialmente en el caso de herramientas que, como las tres analizadas, usan por los poderes públicos como apoyo para tomar decisiones que pueden redundar en restricciones de derechos fundamentales. Y lo segundo, porque también un principio básico del método científico es la publicidad de los métodos, de los datos y de los resultados, para que puedan ser confirmados y refutados por otros investigadores, lo que redundaría en la detección y corrección de errores, la discusión y el debate bien informados y, en último término, la formulación de críticas y propuestas de mejora y el avance del conocimiento. Si una de las razones que se alegan para utilizar algoritmos en apoyo de la toma de decisiones es que la cuantificación y matematización de la información produce decisiones más objetivas y eficaces, y mejores porque están basadas en la evidencia científica, debe ser posible que dicha mayor objetividad y eficacia sean contrastables científicamente, lo que requiere compartir los datos, los códigos y los resultados.

- A este respecto, y de acuerdo con la información disponible, las tres herramientas parecen *a priori* construidas y diseñadas de manera seria y rigurosa. Sin embargo, todas presentan problemas en relación con la transparencia hacia los ciudadanos y hacia la comunidad de expertos, porque faltan datos esenciales para poder realizar una auditoría externa completa totalmente fiable: no se encuentran disponibles todos los *dataset*, no se conoce el diseño del algoritmo ni el peso relativo de cada uno de los factores predictivos y, en el caso de VeriPol, es incluso complicado saber si se usa de manera generalizada y consistente por parte de las autoridades. Dicho esto, es necesario no obstante singularizar el caso de RisCanvi, que merece un juicio aparte y mucho más favorable, en la medida en que la Generalitat de Catalunya está haciendo esfuerzos muy apreciables en los últimos tiempos para suministrar en abierto mucha más información sobre la herramienta y su funcionamiento, como destacamos en el apartado correspondiente.
- Además, la lógica del Reglamento Europeo de Inteligencia Artificial, al que ya hemos hecho referencia, obligaría a una mayor transparencia para el caso de que estas herramientas incorporaran en el futuro IA con capacidad de aprendizaje y adaptativa, y ello en dos sentidos. Por un lado, en cuanto a las obligaciones de una mejor y más completa explicación de la lógica de funcionamiento de las herramientas en abierto con carácter general. Y, por otro lado, y con mucho más detalle, respecto de la información que será obligatorio suministrar a los controladores públicos y autoridades de verificación privadas, para que las administraciones públicas en el ámbito de la Unión Europea puedan emplear sistemas de esta naturaleza.
- En cuanto a la capacidad predictiva y a la eficacia para alcanzar el objetivo propuesto de reducir la comisión de ciertos delitos o la prevalencia de determinadas conductas, es difícil pronunciarse de manera general, porque la cantidad y la calidad de la información disponible sobre cada herramienta analizada son muy dispares, por lo que remitimos al análisis efectuado en el

informe respecto de cada una en concreto. Lo único que destacaríamos aquí como consideración general es que las alegaciones de los responsables del despliegue de estos sistemas, en el sentido de que han mejorado los resultados previos a su introducción, parecen consistentes (con la excepción de VeriPol, donde es imposible cualquier verificación de estos extremos por la falta de información disponible), aunque tienen limitaciones importantes: en algunos casos, la información y los datos sobre los que se sustentan estas afirmaciones son muy escasos, y en general es difícil hacer una comparación con lo que ocurriría en ausencia de estos sistemas. Pero, en todo caso, echamos en falta una reflexión más profunda que compare lo que se gana en términos de prevención del delito con el coste en cuanto a restricción de derechos fundamentales, para poder realizar un juicio sobre la proporcionalidad del empleo de estos sistemas. Juicio que además, en nuestra opinión, debería ser diferenciado según la decisión para la que se utilice el algoritmo: un determinado nivel de acierto en las estimaciones de riesgo puede ser suficiente para tomar decisiones beneficiosas para el condenado o acusado, pero quizá no para tomar decisiones que restrinjan sus derechos; o un determinado nivel de acierto puede bastar para adoptar medidas a escala policial en favor de la víctima, y sin embargo sería discutible que bastara para imponer medidas cautelares al acusado. Consideramos que en este aspecto queda mucho por hacer, y que para que este necesario debate sea siquiera posible es imprescindible un mejor y más completo acceso a la información sobre los algoritmos, su funcionamiento y sus resultados, como hemos señalado *supra*.

- En cuanto a posibles sesgos discriminatorios o impacto dispar de las estimaciones de riesgo sobre grupos diferentes de población, la información publicada por los responsables de estas herramientas es mucho más escasa aún que la disponible sobre capacidad predictiva, y también lo son los datos a partir de los cuales poder hacer una evaluación externa (si bien de nuevo aquí hay importantes diferencias entre las tres herramientas analizadas, siendo RisCanvi –otra vez– la que ofrece más y mejor información; y VeriPol, la más opaca). Nosotros hemos detectado algunos problemas en este sentido, tanto en VioGén como en RisCanvi (y los intuimos en VeriPol, pero aquí la ausencia de datos impide formular nada más que meras intuiciones), si bien las conclusiones que hemos alcanzado han de ser tomadas con mucha cautela, porque están elaboradas sobre una base muy endeble, como hemos tratado de subrayar a lo largo del informe. Ahora bien, precisamente esto es lo que a nuestro juicio resulta más preocupante: que hoy en día (y con la importante excepción parcial de RisCanvi, como destacamos en el apartado correspondiente) resulte imposible estudiar algo tan relevante como los eventuales sesgos o efectos discriminatorios que estas herramientas pueden estar produciendo –o no–, habida cuenta de la ausencia de datos de acceso público con los que trabajar. Que una cuestión de esta relevancia no pueda ser investigada por personas externas a quienes son responsables de la gestión de estas herramientas (y que, por otra parte, ni siquiera ha sido casi estudiada hasta ahora por estos responsables) nos parece un problema grave.

4. PRINCIPALES CONCLUSIONES RESPECTO DE VIOGÉN

Por lo que se refiere a VioGén, las principales conclusiones del estudio son la ya referida necesidad de una mayor información sobre su funcionamiento. Si bien hay varios trabajos sobre distintos aspectos del funcionamiento de la herramienta publicados en revistas académicas de prestigio, eso no permite el tipo de control externo necesario al que nos hemos referido *supra*. En particular, con VioGén podría seguirse un modelo de facilitación de datos anonimizados en abierto similar al que se ofrece, sin problema alguno, con RisCanvi (véase *infra*). O podrían incluirse en el portal estadístico de la web de la Delegación del Gobierno contra la Violencia de Género muchos más datos de los que ahora tiene, que son escasos y aportan poca información. Ello impide una evaluación rigurosa y fiable del funcionamiento de la herramienta desde fuera, tanto en cuanto a su eficacia y fiabilidad como respecto de la identificación de posibles problemas y riesgos.

Con todo, de la información existente se desprende que las afirmaciones de sus responsables sobre el buen funcionamiento de la herramienta en cuanto a su capacidad para reducir la revictimización de las mujeres que denuncian haber sufrido violencia de género parecen consistentes con los datos disponibles en abierto, y con la ampliación de estos que nos ha sido proporcionada. En particular, creemos necesario ser prudentes en relación con las críticas que ha recibido VioGén por no ser capaz de detectar mejor los casos en que la mujer está en riesgo de homicidio a manos de su pareja o expareja. En nuestra opinión, el Ministerio ha realizado un indiscutible esfuerzo por mejorar la predicción de esos casos, que por otra parte son muy difíciles de detectar dada su bajísima prevalencia, e implican siempre una tasa muy elevada de falsos positivos. Sin perjuicio de las mejoras que puedan seguir haciéndose en este ámbito (y que no necesariamente tienen que pasar siempre por la actuación policial o judicial), nos parece necesario ser conscientes de los límites a que la predicción de este tipo de eventos está sometida, y reconocer que ningún sistema de estimación de riesgos puede aspirar razonablemente a reducir estas tragedias a cero.

La capacidad predictiva de VioGén es buena por lo que hace a la detección del riesgo bajo (en el sentido de que en la inmensa mayoría de los casos de riesgo bajo no hay agresiones posteriores), pero más débil por lo que se refiere a los casos de riesgo alto, pues muchos casos en los que se afirma riesgo alto o extremo no van seguidos tampoco de agresiones ulteriores. Sin perjuicio de reconocer que la propia implementación de medidas de protección contribuye seguramente en parte a este efecto, hay que ser conscientes de que el sistema, como reconocen sus propios responsables, asume un nivel elevado de falsos positivos para poder maximizar la sensibilidad (esto es, no dejar sin protección a ninguna mujer que realmente pueda estar en riesgo). Y esta opción, que policialmente no es problemática desde la perspectiva de los derechos fundamentales, porque las medidas de protección de las mujeres inciden muy poco en los derechos de los varones denunciados, resultaría sin embargo mucho más cuestionable si se utilizaran los niveles de riesgo de manera automática en el ámbito judicial para fundamentar la imposición de medidas cautelares, especialmente las más restrictivas de derechos. Esta utilización generalizada no está ocurriendo hoy en día, al menos hasta donde hemos podido

averiguar, pero es algo frente a lo que conviene estar alerta porque, en primer lugar, ya en la actualidad es posible tener en cuenta el resultado de VioGén para tomar estas decisiones, aunque no se perciba como determinante; y, en segundo lugar, porque es probable que este uso pase a ser mayor en el futuro. En tal caso, el juicio emitido en este sentido debiera ser, como es obvio, notablemente revisado.

Por lo que hace a posibles sesgos o efectos discriminatorios de VioGén, ya hemos dicho que la muy limitada información disponible dificulta mucho poder analizar esta cuestión con un mínimo de rigor. Nosotros hemos detectado indicios que apuntan a una protección menor, o menos eficiente, de las mujeres nacidas fuera de España, tanto frente a nuevas agresiones de género como frente a homicidios. Señalamos en el estudio la metodología empleada para analizar los datos disponibles y las razones por las que llegamos a esta conclusión, que ciertamente se basa en información incompleta, pues incompleta es la información disponible al respecto, pero que consideramos que, en ausencia de más datos, y más completos, debería llevar al menos a una revisión del funcionamiento de la herramienta, tanto para la estimación del riesgo de violencia en general como para la estimación del riesgo de homicidio con la nueva escala H.

En cuanto al cumplimiento normativo, se echan en falta evaluaciones más completas de riesgos, tanto en cuanto a los datos de carácter personal empleados y los riesgos que de su uso puedan derivarse como, sobre todo, en materia de las evaluaciones sobre riesgos para derechos fundamentales y de prevención de sesgos que un sistema de este tipo debería incorporar. De hecho, en el modelo de cumplimiento del Reglamento Europeo en materia de inteligencia artificial estas evaluaciones pasarán a ser obligatorias para sistemas de esta índole desde el momento en que incorporen IA con capacidad de evolución o aprendizaje (lo que la versión actual de VioGén parece que no hace, pero que en cuanto pase a ser empleada una que sí lo haga deberá adaptarse a esta exigencia). Por último, también son susceptibles de una mejora importante los mecanismos que garanticen el acceso efectivo de los ciudadanos a sus evaluaciones y a todos sus datos de carácter personal almacenados por el sistema, algo que en la actualidad no se garantiza suficientemente.

5. PRINCIPALES CONCLUSIONES RESPECTO DE RISCANVI

RisCanvi es la herramienta más transparente de las analizadas. Si bien el diseño de los algoritmos se mantiene secreto (con la excepción del que estima el riesgo de reincidencia violenta, que hemos podido conocer gracias a la colaboración de colegas de la Universitat Pompeu Fabra), la Generalitat de Cataluña ha ido poniendo en práctica, de manera creciente con los años, una política de *open data* con publicación de datos anonimizados en formato Excel en su página web, que es reseñable en un sentido muy positivo, que creemos que debería servir de modelo a otras administraciones, y que permite analizar el funcionamiento de los algoritmos a cualquier investigador interesado sin comprometer en modo alguno la privacidad de las personas afectadas por las estimaciones de riesgo. Además, la herramienta ha sido auditada en cuanto a su funcionamiento recientemente, y los resultados de esta evaluación han sido también publicados. Desde este punto de vista, si en un futuro la

herramienta incorporara para la mejora de sus estimaciones de riesgo sistemas de IA en el sentido definido por el Reglamento europeo, probablemente podrá cumplir las exigencias de transparencia y auditabilidad del sistema con bastante facilidad.

Queremos subrayar que de los cinco tipos de riesgo que estima RisCanvi, en este informe solo hemos analizado el de reincidencia violenta, por lo que todas nuestras consideraciones se limitan exclusivamente a este aspecto.

La capacidad predictiva de RisCanvi es bastante buena en cuanto a la detección de las personas con riesgo bajo, pero mucho menor en cuanto a la detección de las personas con riesgo alto. Las estimaciones de riesgo incluyen un número elevado de falsos positivos y en esto coinciden tanto las investigaciones llevadas a cabo desde organismos dependientes de la Generalitat, como los cálculos que nosotros hemos podido realizar a partir de los datos disponibles, así como los resultados de la auditoría externa. Creemos que los responsables de la herramienta son conscientes de esta realidad y están extrayendo de ella la conclusión, razonable desde nuestro punto de vista, de que las evaluaciones de riesgo deberían conducir a una suavización de los regímenes de vida y a mayores posibilidades de acceso a la libertad para las personas de riesgo bajo, pero no llevar automáticamente a un endurecimiento del régimen ni de las condiciones de cumplimiento para aquellos clasificados como de riesgo alto. Sería deseable, no obstante, que existieran pautas claras de actuación a este respecto, y también que este conocimiento se extendiera al ámbito judicial, donde se revisan las decisiones de clasificación de las juntas de tratamiento, y donde la Fiscalía tiende a dar una importancia muy elevada a las clasificaciones de riesgo alto o medio, sin tener en cuenta las importantes limitaciones que esas estimaciones tienen.

Algunos de los factores que RisCanvi tiene en cuenta para evaluar el riesgo afectan muy directamente a la vida privada e intimidad de las personas y a su libre desarrollo de la personalidad o a su entorno de socialización. Puesto que no son públicos los algoritmos (con la excepción antes mencionada) no se puede saber cuál es la capacidad predictiva de esta información tan sensible. Y sin embargo puede resultar cuestionable que aspectos sobre los que la persona no tiene control (antecedentes delictivos de sus familiares, por ejemplo) o forman parte de su derecho al libre desarrollo de la personalidad (número y tipo de actividades sexuales preferidas, por ejemplo) puedan constituir una base legítima para fundamentar medidas restrictivas de derechos (por ejemplo, la denegación del acceso a un tercer grado o a la libertad condicional). Sería necesario comparar el rendimiento de estos ítems (es decir, cuánto se pierde en cuanto a capacidad predictiva, si se eliminan del algoritmo) con su coste en términos de afectación a los derechos (a la intimidad, a la igualdad, al respecto a la vida privada y familiar, etc.), para tomar una decisión sobre si resulta proporcionado emplearlos en las predicciones. Además, en algunos de los ítems puede incluso cuestionarse que sea pertinente tenerlos en cuenta incluso en el supuesto de que tengan valor predictivo (por ejemplo, la presencia de discapacidad). Creemos que sería necesaria una auditoría profunda respecto del riesgo que comportan estos elementos para los derechos fundamentales de los reclusos, que no nos consta que se haya llevado a cabo, así como una transparencia mucho mayor sobre las evidencias empíricas que, en su caso, podrían avalar que se entiendan estos factores como estadísticamente significativos.

Por último, nuestro análisis estadístico ha revelado cierto trato extraordinariamente severo para algunos colectivos de reclusos, que no parece justificado y que podría retroalimentar sesgos en la herramienta en perjuicio de estos colectivos, esencialmente las personas con problemas de salud mental o de adicciones, o personas con una situación socioeconómica muy desfavorable. No es el caso de Ris-Canvi, y creemos que esto es una virtud propia que debe ser destacada, en cuanto a un trato más severo con los extranjeros, que no se da según los datos disponibles. Entendemos que el funcionamiento de la herramienta debería al menos revisarse a partir de los datos estadísticos que aportamos, y valorarse alguna vía para solucionar estos problemas.

6. PRINCIPALES CONCLUSIONES RESPECTO DE VERIPOL

VeriPol es la herramienta de la que menos información se dispone, lo que agrava el problema de falta de transparencia y las dificultades, cuando no imposibilidad, para realizar la tarea de control externo que hemos pretendido con este informe. Agradecemos a la Policía Nacional la información y los datos que nos proporcionó para al menos podernos hacer alguna idea sobre esta herramienta y el tipo de uso que se le da en la práctica, pero aun así se trata de una situación muy deficiente en cuanto al cumplimiento de los estándares mínimos de transparencia que consideramos requeridos para el empleo de herramientas de este tipo, y aunque en el informe hemos analizado toda la información de la que hemos podido disponer, muchas conclusiones no pasan de ser conjeturas, a falta de datos contrastables con los que poder trabajar.

En cuanto a los problemas que pueden derivarse de su afección a derechos fundamentales o en cuanto a la posible generación de sesgos, esta misma falta de información impide una valoración externa sobre el funcionamiento práctico del sistema. El informe ha tratado, al menos, de establecer un marco teórico respecto de las cautelas mínimas exigibles para evitar que modelos como este, basados en IA y en aprendizaje a partir de sistemas de procesamiento de lenguaje natural, incurran en sesgos para las funciones que se le presuponen (detección de denuncias falsas por parte de los ciudadanos y ciudadanas), que, por ejemplo, puedan afectar más a la parte de la población que maneja peor el idioma por razones que pueden ser de muy diverso tipo (nivel socioeconómico, origen, nivel de estudios e incluso la mera diferencia diatópica dialectal en aquellos casos en que pueda emplearse una variedad diferente a la que ha sido mayoritariamente empleada para entrenar el modelo). Asimismo, y como para el resto de las herramientas, trata de sugerir mejoras para una mejor evaluación de los posibles riesgos para los derechos fundamentales que puedan derivarse de su utilización, riesgos tanto mayores cuanto más extendido sea su uso.



¿Cuál es la capacidad predictiva de los algoritmos diseñados para prevenir la delincuencia? ¿Son eficaces? ¿Son lo suficientemente transparentes? ¿Tienen sesgos? Para responder a estas preguntas, este estudio analiza tres herramientas algorítmicas utilizadas por la Administración española para hacer estimaciones de riesgo en los ámbitos policial y penitenciario (VioGén, RisCanvi y VeriPol, esta última con inteligencia artificial) con la finalidad de evaluar su funcionamiento, en particular en términos de su respeto a los derechos fundamentales.

El estudio combina métodos cuantitativos y cualitativos, y ha sido realizado por profesores de la Universitat de València y de la Universitat Politècnica de València, expertos en distintas ramas del derecho (penal, administrativo y constitucional), así como en estadística y matemáticas. Los objetivos principales son analizar la capacidad predictiva de los tres sistemas y su eficacia para alcanzar el objetivo que pretenden (reducción de la delincuencia); valorar su grado de transparencia, tanto para la ciudadanía en general como para las personas directamente afectadas por su uso; identificar y discutir posibles sesgos, y comprobar si respetan el marco legal europeo y español vigente, prestando especial atención a la normativa sobre protección de datos personales y a la regulación del uso de sistemas automatizados de toma de decisiones.

También se analiza si los costes para los derechos fundamentales que conlleva el uso de estas herramientas exceden lo que sería legítimo, necesario y/o proporcionado en relación con el beneficio actual o probable en términos de prevención del delito. El estudio subraya la importancia de la transparencia para hacer posible el necesario debate público que debe existir sobre herramientas que, como las tres analizadas, son utilizadas por los poderes públicos como apoyo para tomar decisiones que pueden redundar en restricciones de derechos fundamentales.