

# Un modello di *machine learning* per l'identificazione di aziende collegate alla criminalità organizzata in Italia, sulla base di dati di bilancio

SINTESI NON TECNICA DELLO STUDIO

“A MACHINE LEARNING APPROACH FOR THE DETECTION OF FIRMS LINKED TO ORGANISED CRIME IN ITALY, BASED ON BALANCE SHEET DATA”

P. Cariello, M. De Simoni, S. Iezzi\*

Lo studio ha l'obiettivo di sviluppare un algoritmo di apprendimento automatico (*machine learning*) per individuare le aziende italiane potenzialmente connesse a contesti di criminalità organizzata.

Il modello è rivolto a imprese con la forma giuridica di società di capitale. Per l'addestramento viene utilizzato un dataset di 1.804.278 imprese italiane, relativo al periodo 2010–2021, che incorpora dati di bilancio, informazioni sulla esposizione debitoria nei confronti del sistema bancario e finanziario, dati occupazionali e sull'assetto proprietario e di *governance*. Per l'addestramento dell'algoritmo è stato utilizzato un campione di 28.570 imprese ad alto rischio di collegamento con la criminalità organizzata; la selezione del campione è basata su fonti pubbliche e, in larga misura, su una selezione della cd. 'mappatura' delle imprese potenzialmente connesse a contesti di criminalità organizzata, sviluppata presso l'UIF (cfr. Rapporto Annuale UIF sul 2020, pagg. 47-48). Questo campione di addestramento rappresenta uno dei più ampi censimenti di imprese di questo tipo, e costituisce una delle innovazioni più rilevanti dello studio<sup>1</sup>.

Come avviene in approcci di tipo *machine learning*, le prestazioni dell'algoritmo utilizzato per l'addestramento (XGBoost) sono state calibrate con l'obiettivo di massimizzare uno dei parametri che misurano la capacità previsiva del modello; in questo caso si è scelto il cd. tasso di *recall* (o *sensitività*), che misura l'incidenza di imprese riconosciute dal modello come infiltrate sul totale delle imprese effettivamente infiltrate prese in esame. Utilizzando come campione di test l'insieme di imprese infiltrate non utilizzate per l'addestramento (cd. test *out-of-sample*<sup>2</sup>), si ottiene un tasso di *recall* pari al 75,6%: in altri

---

\* Unità di Informazione Finanziaria per l'Italia.

<sup>1</sup> L'infiltrazione o collegamento alla criminalità organizzata possono essere accertati solo a livello investigativo e giudiziario; tuttavia, per semplicità di esposizione, nel prosieguo di questa sintesi ci riferiremo al campione di addestramento come campione di imprese 'infiltrate'.

<sup>2</sup> Per testare la validità del modello, il campione iniziale di imprese infiltrate viene diviso in due sottoinsiemi: circa 260 mila osservazioni per impresa/anno sono usate per addestrare il modello, le restanti 66 mila per verificare la sua capacità di riconoscere le imprese infiltrate. Questa operazione (divisione del campione in due sottoinsiemi, addestramento del modello su una parte delle imprese e verifica della sua validità sulle restanti imprese) è stata effettuata cinque volte.

termini, su un totale di 32.166 osservazioni per impresa/anno di imprese infiltrate (distribuite nel periodo 2010-2021), il modello è stato in grado di riconoscerne correttamente 24.309. Un altro importante parametro è la *specificity*, ossia la quota di imprese non-infiltrate riconosciute come tali dal modello, che risulta pari al 74,2% (ossia, su un totale di 33,410 osservazioni di imprese non-infiltrate, il modello ne riconosce correttamente 24.801). Questi risultati non si discostano molto da quelli ottenuti inizialmente addestrando il modello sull'intero campione di imprese infiltrate (cd. test *in-sample*), fornendo segnali incoraggianti sulla stabilità del modello e sulla sua capacità di mantenere buoni livelli di prestazione anche se applicato a nuovi dati.

L'algoritmo consente di calcolare un indicatore di rischio per oltre 900 mila società di capitale attivo in Italia (sono stati usati i dati dei bilanci 2021): tale *score* è compreso tra 0 a 1 e può essere interpretato, indicativamente, come la probabilità che la singola impresa sia connessa a contesti di criminalità organizzata. Il 78,3% delle società di capitale ha uno score inferiore a 0,50; tra il restante 21,7% delle aziende con un punteggio superiore a tale soglia, l'1,8% ha uno score superiore a 0,95.

Sono stati anche effettuati diversi esercizi di validazione del modello con dati indipendenti. In particolare, i risultati delle prime versioni statistiche del modello sono stati confrontati massivamente sia con i dati UIF delle Segnalazioni di operazioni sospette, sia con le evidenze del Nucleo Speciale di Polizia Valutaria della Guardia di Finanza, con esiti sostanzialmente positivi. Da ultimo, la versione più aggiornata del modello è stata validata utilizzando alcuni dati relativi a 1) le aziende colpite da provvedimenti prefettizi di interdittiva antimafia e, all'opposto, 2) quelle incluse nelle cd. "*white list*" (ossia imprese operanti in specifici settori economici, maggiormente esposti al rischio di infiltrazione mafiosa, per le quali è attestata a livello prefettizio l'assenza di connessioni note con la criminalità organizzata). Le imprese soggette a interdittiva presentano uno *score* di rischio mediano che supera di oltre due volte e mezzo quello delle imprese incluse nelle "*white list*". Inoltre, il modello individua il 64,6% delle imprese soggette a interdittiva come infiltrate; simmetricamente, il 70,5% delle imprese nelle "*white list*" è riconosciuto come non infiltrato dal modello. L'ordine di grandezza di questi riscontri, ottenuti 'sul campo' con dati del tutto indipendenti dallo studio, non si discosta troppo da quello delle misure *out-of-sample* di *recall* e *specificity* prima riportate.

L'indicatore di rischio proposto nello studio – che è ancora in versione sperimentale - ha varie potenziali applicazioni. In ambito strategico, può consentire ad esempio l'elaborazione di mappe di rischio a livello territoriale o settoriale. In ambito operativo, può contribuire al patrimonio informativo che supporta le funzioni istituzionali dell'UIF; potrebbe anche essere utilizzato come strumento preliminare di *screening* per contribuire a orientare l'azione degli organi investigativi, ad esempio nel monitoraggio dell'utilizzo dei fondi pubblici (PNRR). Conferme della sua validità operativa dovranno tuttavia venire da ulteriori applicazioni 'sul campo'.