

# Data Privacy for AI Fraud Detection Models

## A framework for GDPR compliant AI

Kadir Ider

Delivery Hero SE

Berlin, Germany

email: kadir.ider@deliveryhero.com

Andreas Schmietendorf

Berlin School of Economics and Law

Berlin, Germany

email: andreas.schmietendorf@hwr-berlin.de

**Abstract**—Although the European General Data Protection Regulation (GDPR) is technology neutral, it indirectly imposes strict processing rules for personal information in Artificial Intelligence systems. As fraud detection becomes more sophisticated and complex, the challenge to manage the trade-off between privacy and accuracy of such systems arises concurrently. This paper identifies and presents key components for a GDPR compliant design and development of Machine Learning supported fraud detection solutions.

**Keywords** - *Fraud detection; GDPR; transparency; automated decision making; AI.*

### I. INTRODUCTION

According to a global economic crime and fraud survey in 2020, the financial damage in the US is estimated to amount to USD 42 billion, with a large fraction contributed by cybercrime [13]. Nearly half of businesses have been affected by fraud attacks within a 24 months period. At the same time false credit card declines amount to almost USD 120 billion that is three times higher than the detection of actual fraud cases [9].

Thus, there is a great interest in increasing the accurate detection of true positive and decreasing false positive fraud cases to achieve both, minimize monetary losses, but also accelerate business. Roughly 50% of affected corporations employ Artificial Intelligence (AI) fraud detection systems, but struggle to harness the benefits of such tools. In order to fight fraud effectively, algorithms must be provided with adequate input data that include personal identifiers to some extent.

The European Data Protection authorities are aware of this trend and increasingly publish guidelines on the lawful usage of personal Identifiable Information (PII) as well as to prevent the “black box problem” [12]. The phenomenon of black box algorithms is due to increasing sophistication and complexity of Machine Learning (ML) solutions making it difficult to transparently process PII and ensure accountability [15].

Particularly when sensitive PII pursuant to Art. 9 GDPR [16] is processed, the data processing may not be valid if the legal basis is not given and if it is not consistent with the initial purpose of the data collection. Compliant fraud detection systems can therefore act as an instrument to meet

accountability requirements by identifying unlawful and non-compliant usage of PII.

The body of the paper is divided into six sections. In Sections 2 and 3, a segmentation of fraud activities and models generate structural insights and define the research scope. Section 4 evaluates relevant papers and identifies the research need for the framework development. Section 5 formulates the framework components, based on the underlying requirements. In Section 6, further specification is provided on the framework applicability. The final section highlights and discusses the trade-off of privacy requirements and usage of AI and ultimately closes with the conclusion.

### II. FRAUD ACTIVITIES AND METHODS CLASSIFICATION

Frequent online fraud activities include identity theft, account takeovers, abuse of promotions, fake reviews or listings [2]. All have monetary consequences in common, i.e., financial losses due to fraud and potential fines imposed to the data controlling entity as a consequence of the data theft. Fraud detection methods include

- 1) *blacklists,*
- 2) *rule engines and*
- 3) *AI solutions.*

Blacklists could contain user data, such as name, email, IP address and device data that are associated with fraudulent activities and therefore will be blocked from using online services. In the scope of promotion code abuse, such lists may be effective if users have already created an account and where one or more variables are matching the blacklist. Due to its static characteristic, the list may only be effective after a fraud attempt has been made, thus may have already caused damage at the stage of detection. Such a reactive approach is therefore not sufficient as a standalone solution. Blacklists may be considered as the most simple rule engine.

More complex, but also manually written rule engines employ several rules, check multiple conditions and incorporate weight scorings. Such rules are frequently employed in the detection of fraudulent activities in the scope of money transactions. For example, multiple small instead of large transactions conducted by different people from an unusual location or having the same beneficiary would be probably rejected as the transaction would violate one or many rules [11]. The downside is that the larger the number of rules, the greater the maintenance. Moreover,

rules may cancel each other out. Both options are suitable for identifying obvious fraud cases, are computationally cheaper, but require more manual work to maintain their effectiveness. On the upper end of proactive fraud detection approaches are Supervised Machine Learning algorithms (SML), such as Decision Trees, Support Vector Machine (SVM) or Artificial Neural Networks (ANN) [11]. To overcome the burden of maintaining rule engines, SML models learn from existing patterns and identify fraud in unstructured data, learn and predict fraud activities, despite a multitude of input features [1].

### III. SCOPE AND DEFINITIONS

In order to limit the research scope and ensure thorough analysis, the assessment will focus on the development of an AI compliance framework, but will not further elaborate the design of a comprehensive Data Protection Management System (DPMS). Subsequently, according to Figure 1, the third-party management will not be assessed further as it belongs to the DPMS framework.

No matter which fraud detection algorithm is used, all are gaining popularity due to their ability to exploit large amounts of personal data, conduct automated decision making and create profiles. Such activities provide competitive advantage and leverage business activities. However, these kinds of processing activities require additional security measures to protect PII, pursuant to Art. 22 and 35 GDPR [16]. As not only data of identified, but also identifiable persons are affected, the GDPR sets strict requirements on such activities. For example, in the design and development process of fraud detection technology, data protection by design and by default, in accordance with Art. 25 GDPR [16], requires enterprises to design the solution in such a way that the flow of personal data is protected by Technical and Organizational Measures (TOMs) at any point in time during the processing, as stated in Table 1. The integration of cloud solutions - whether internally or externally developed and/or hosted - must additionally protect data in rest and transition.

In order to consistently demonstrate compliance and accountability in the entire data lifecycle, monitoring and incident response management plans for the detection of fraudulent activities must be documented, even after implementation of technologies. Non-compliance may result in data breaches that could compromise the identity of data subjects and lead to increased risks to the rights and freedoms of individuals, pursuant to Recital 75 GDPR. Monetary penalties could arise for the data controlling and processing enterprise and amount to €20 million or 4% of the total worldwide annual turnover of the preceding financial year, whichever is higher, according to Art. 83 (5) GDPR [16].

The British Information Commissioner's Office (ICO) has published guidelines and frameworks for AI audits which will be used to a large extent for the development of a comprehensive framework, including specifications for the proposed control areas.

### IV. LITERATURE REVIEW

In related publications, either the development, assessment of features, comparison of fraud detection algorithms performance or the issue of an increased digital footprint and its data protection implications are elaborated. A comparison of privacy preserving fraud detection methods [1] provides an understanding on the performance of algorithms by measuring their efficiency. GDPR or data protection aspects are marginalized. Other publications assess elements of trustworthiness in the usage of AI or look into specific AI models including respective classification techniques, such as for ANN, aiming to improve their prediction accuracy [7][10]. The mentioned articles do not highlight the importance of transparency and accountability for PII adherent to the GDPR. Thus, the fraud detection algorithms proposed in the articles disregard the fact that the usage of such technology introduces new risks to data, but more importantly to individuals, which leads to the development of a framework to close the present shortcomings. The development shall further improve the accountability requirements pursuant to Art. 5 (2) GDPR [16].

### V. ASSESSMENT METHOD

In the following assessment process, a GDPR compliant AI audit framework will be presented. Subsequently, a common data lifecycle process is analyzed independent of a particular fraud detection model as the framework is equally applicable. The ICO's proposal for a compliant AI framework includes an assessment of the general governance and accountability aspects and AI specific control areas [14]. The model's drawback is the lack of an ethical control area, which will be compensated by incorporating key requirements into the framework based on the Hambach Declaration on Artificial Intelligence [5].

The proposed AI Privacy Design Framework enhances an already existing DPMS [17]. The AI component is split into nine elements and looks as depicted in Figure 1.

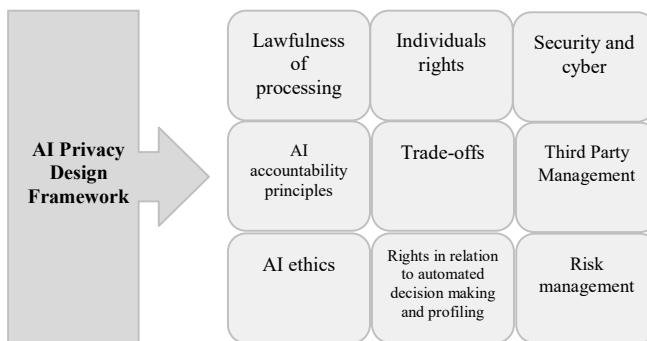


Figure 1. AI Privacy Design Elements. Source: own elaboration.

Each element is a standalone feature, which increases the compliance with GDPR, when considered in the framework. However, removing an element immediately reduces the quality of the framework's proficiency, overall. The composed elements and corresponding specifications are

non-exhaustive and should be adjusted according to the business needs. The elements displayed in Figure 1 are further specified in Table 1 (see Appendix 1). The framework can be effectively used for the definition of guidelines or development of maturity assessment models.

## VI. PRIVACY PRESERVING DESIGN ASPECTS

The three parties that play an essential role in the process are data subjects, data controller and data processor. The data controller pursuant to Art. 4 (7) GDPR [16] is the entity that defines the purpose and means of processing and is liable for any data incident, while the processor is another party that processes the data upon instructions of the controller. Typically, the processor is an entity that either hosts the data in large data centers (on behalf of the controller) or provides cloud services, such as - pretrained - fraud detection models.

That concludes that the controller is in the obligation to design and ensure a privacy preserving procedure throughout the PII usage lifecycle. In accordance with the framework presented in Table 1, the lifecycle begins with the determination and definition of the lawful basis associated with the data subject rights on the subsequent data processing. Most relevant lawful basis in the scope of fraud detection are legitimate interest and data subjects consent. Processing PII for training and testing based on legitimate interest provides the controller the broadest legitimate ground. It allows the utilization of PII to its full extent for testing various fraud detection purposes, their prediction accuracy as well as use the data for a wider range of AI-based models. Yet, because of its flexibility, it may not be the most accurate ground for processing.

A three-stage-test should be performed to test the fitness of this legitimate basis and involves [8]:

- 1) *identify a legitimate interest (the ‘purpose test’);*
- 2) *show that the processing is necessary to achieve it (the ‘necessity test’); and*
- 3) *balance it against the individuals interests, rights and freedoms (the ‘balancing test’).*

If the conclusion favours the interests of the controller, the legitimate interest may be appropriate.

Reliance on consent is appropriate in cases where the deployment of fraud detection is in the immediate context with the data subject, e.g., the prevention of customer account or credit card misuse. Processing PII requires the collection of separate consent as each activity has a different purpose and might require the processing of different PII.

Pursuant to Art. 7 GDPR [16], the conditions for obtaining valid consent are:

- 1) *freely given,*
- 2) *specific,*
- 3) *informed and unambiguous,*
- 4) *clear affirmative act of the individual (e.g., clicking “I consent”).*

The downside is that the number of different purposes increases the difficulty to ensure the conditions of Art. 7 GDPR [16] are effectively met. The data subject has the right to restrict the processing or withdraw the consent completely at any time. Consequently, the immediate discontinuation of

the data processing for the fraud detection purpose must be ensured (see Table 1 “Individuals rights”).

Besides determining the lawful basis, the controller must challenge whether or not the intended data is needed for the development and deployment of the fraud detection model in accordance with the minimization principle presented in the section “AI accountability principles” in Table 1. Common data categories selected in the feature engineering process include identity, orders, payment method, location, network data and [2]. These parameters are all considered as PII, as all data are in association with an individual.

Acquiring fraud detection services from external suppliers does not release the controller from the duty to adhere to GDPR compliance. The obligation involves defining and communicating the requirements down to the processor.

In order to strengthen the protection of PII, it is suggested to implement additional security mechanisms, such as homomorphic encryption [6]. The original dataset will be encrypted, but still provides the ability (for the algorithm) to perform computations on the encrypted data. However, depending on the homomorphic encryption method, i.e., full, somewhat or partial, the computational overhead may slow down the entire fraud detection process [3]. The optimal trade-off between model complexity and PII security must therefore be balanced out. Nevertheless, encryption is considered as a pseudonymization of PII, as a re-identification is possible with the corresponding decryption key, at any time.

## VII. ADVANTAGES AND LIMITS OF AI SET BY PRIVACY REGULATIONS

Compliance with GDPR is at the ultimate forefront of advantages. The chapters highlight that, by adhering to the data protection law, AI deploying enterprises demonstrate thorough understanding of their models. A transparent documentation of fraud detection models will not reveal secrets about the underlying algorithm and thus, will not jeopardize businesses intellectual property, but rather enable affected users to understand the processing of personal information. From a business continuity point of view, a compliance and thus transparent documentation increases the maintainability of fraud detection models, particularly in areas with higher employee turnover rates.

A clear limitation of AI models is set by the consistent challenge to measure the trade-offs between the level of privacy and model accuracy. Implementing security layers, such as homomorphic encryption elevate the data security, but come at the expense of speed and required computation resources. These factors may decrease the appeal for smaller and medium enterprises with fewer tech- and privacy experts as well as financial resources. Companies that deploy application programming interface (API) based fraud detection models must additionally monitor user queries, implement rate-limiting and other security layers [8]. Moreover, deploying externally developed models binds to their bias and therefore might not effectively detect the fraud

as the degree of influence in the development is limited and set by the provider.

### VIII. CONCLUSION

The foregoing analysis presented a framework for a GDPR compliant development of fraud detection models. Incorporating privacy and ethics into technology goes beyond the mere understanding of the 99 GDPR articles. The challenge is - irrespective of the underlying Machine Learning model - to translate the regulation requirements into operable measures. Due to the GDPR's technology neutrality, enterprises face the great challenge to identify and integrate appropriate technology in order to comply. The proposed framework must therefore not be seen as a standalone solution. The eight control areas defined are rather modules to be incorporated in the scope of AI development and where PII is involved. In subsequent studies, the implementation and assessment of this framework on various fraud detections, but also other AI supported models in different industries will eventually reveal its long-term effectiveness. In this context, an empirical validation of the proposed framework elements will be further conducted.

### REFERENCES

- [1] M. P. Bach, N. Vlahović, and J. Pivar, "Fraud Prevention in the Leasing Industry Using the Kohonen Self-Organising Maps." *Organizacija* 53, no. 2, May 2020, pp. 128–45. <https://doi.org/10.2478/orga-2020-0009>.
- [2] M. Beals, M. DeLiema, and M. Deevy, "Full-Taxonomy-Report" [retrieved: 10, 2020]. <http://162.144.124.243/~longevl0/wp-content/uploads/2016/03/Full-Taxonomy-report.pdf>.
- [3] R. Canillas et al., "Exploratory Study of Privacy Preserving Fraud Detection." In *Proceedings of the 19th International Middleware Conference Industry on - Middleware '18*, pp. 25–31. Rennes, France: ACM Press, 2018. <https://doi.org/10.1145/3284028.3284032>.
- [4] Data2.eu, "What Technical and Organisational Measures Do We Need to Take? - GDPR Processing Index." [retrieved: 10, 2020]. <https://data2.eu/en/gdpr/what-technical-and-organisational-measures-do-we-need-to-take>.
- [5] Datenschutzkonferenz (DSK), "Report on Experience Gained in the Implementation of the GDPR." [retrieved: 10, 2020]. [https://www.datenschutzkonferenz-online.de/media/dskb/20191213\\_evaluation\\_report\\_german\\_dpa\\_s\\_clean.pdf](https://www.datenschutzkonferenz-online.de/media/dskb/20191213_evaluation_report_german_dpa_s_clean.pdf).
- [6] T. Graepel, K. Lauter, and M. Naehrig, "ML Confidential: Machine Learning on Encrypted Data", In *Information Security and Cryptology – ICISC 2012*, edited by T. Kwon, M.-K. Lee, and D. Kwon, 1–21. Berlin, Heidelberg: Springer Berlin Heidelberg, 2013.
- [7] C. W. L. Ho, J. Ali, and K. Caals, "Ensuring Trustworthy Use of Artificial Intelligence and Big Data Analytics in Health Insurance." *Bulletin of the World Health Organization* 98, no. 4, Apr. 2020, pp. 263–69. <https://doi.org/10.2471/BLT.19.234732>.
- [8] Information Commissioner's Office, "Guidance on the AI auditing framework." [retrieved: 10, 2020]. <https://ico.org.uk/media/about-the-ico/consultations/2617219/guidance-on-the-ai-auditing-framework-draft-for-consultation.pdf>.
- [9] J. Khodos, "Mastercard Rolls Out Artificial Intelligence Across Its Network." [retrieved: 10, 2020]. <https://newsroom.mastercard.com/press-releases/mastercard-rolls-out-artificial-intelligence-across-its-global-network/>.
- [10] A. K. Patience and J. V. N Lakshmi, "A Study on Credit Card Fraud Detection Using Machine Learning." *International Journal of Trend in Scientific Research and Development (Ijtsrd)* 4, no. 3, Apr. 2020, pp. 801–804.
- [11] M. M. Rahman and A. R. Saha, "A Comparative Study and Performance Analysis of ATM Card Fraud Detection Techniques." *Journal of Information Security* 10, no. 03, Jul. 2019, pp. 188–97. <https://doi.org/10.4236/jis.2019.103011>.
- [12] A. Rai, "Explainable AI: From Black Box to Glass Box." *Journal of the Academy of Marketing Science* 48, no. 1, Jan. 2020, pp. 137–41. <https://doi.org/10.1007/s11747-019-00710-5>.
- [13] K. Rivera, C. Rohn, J. Donker, and C. Butter, "PwC's Global Economic Crime and Fraud Survey 2020," 2020. [retrieved: 10, 2020]. <https://www.pwc.com/gx/en/forensics/gecs-2020/pdf/global-economic-crime-and-fraud-survey-2020.pdf>
- [14] B. Reuben and V. Gallo, "An Overview of the Auditing Framework for Artificial Intelligence and Its Core Components." [retrieved: 10, 2020]. <https://ico.org.uk/about-the-ico/news-and-events/ai-blog-an-overview-of-the-auditing-framework-for-artificial-intelligence-and-its-core-components/>.
- [15] M. Veale, R. Binns, and L. Edwards, "Governing Artificial Intelligence: Ethical, Legal and Technical Opportunities and Challenges." *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences* 376, no. 2133, Nov. 2018, 20180080. <https://doi.org/10.1098/rsta.2018.0080>.
- [16] EU General Data Protection Regulation (GDPR): Regulation (EU) 2016/679 of the European Parliament and of the Council of 27 April 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing Directive 95/46/EC (General Data Protection Regulation), OJ 2016 L 119/1.
- [17] A. Rai, "Explainable AI: From Black Box to Glass Box." *Journal of the Academy of Marketing Science* 48, no. 1, Jan. 2020, pp. 137–41. <https://doi.org/10.1007/s11747-019-00710-5>.

APPENDIX 1

TABLE I. AI PRIVACY DESIGN FRAMEWORK

<i>Areas and controls</i>	<i>Specification</i>
<b>Lawfulness of processing</b>	Pursuant to Art. 6 GDPR [16]
Assessment of lawful basis	<ul style="list-style-type: none"> <li>• Determine lawful basis depending on the fraud systems purpose:                             <ul style="list-style-type: none"> <li>○ Determination before processing starts</li> <li>○ Document the decision process</li> <li>○ Lawful basis must not be changed after processing starts</li> <li>○ Communicate lawful basis with affected individuals (e.g., via privacy policy)</li> </ul> </li> </ul>
<b>AI accountability principles</b>	Pursuant to Art. 5, 13, 14 and Recital 60 GDPR [16]
Fairness and transparency in profiling	<ul style="list-style-type: none"> <li>• Specification of person/team in charge of AI system (responsibility and accountability)</li> <li>• Explanation of models unbiased decision making ability</li> <li>• Address (e.g., in privacy policy) the associated risks of using AI                             <ul style="list-style-type: none"> <li>○ Providing a summary of a data protection impact assessment (DPIA)</li> </ul> </li> </ul>
Accuracy (of used data)	<ul style="list-style-type: none"> <li>• Specification of data used in AI system (classification in categories and detailed listing of dates, e.g., Identification Data: name, surname; Technical Data: IP address, device ID)</li> <li>• “Binning” (e.g., continuous) variables into discrete ranges in the pre-processing phase may alter the accuracy of data (pursuant to Art. 5 (1) lit. d GDPR [16] and subsequently the accuracy of the prediction (e.g., instead of processing individuals age “54”, he/she will be “binned” into the age group “50-60”))</li> </ul>
Data minimization and purpose limitation	<ul style="list-style-type: none"> <li>• Ability to explain why AI is required for the specific purpose (e.g., “detect the abuse of promo codes”, instead of just stating “fraud detection”)</li> <li>• Aiming for usage of minimum amount or anonymous data, if sufficient enough for achieving a specific legitimate purpose</li> </ul>
<b>AI ethics [5]</b>	Pursuant to Art. 5, 12, 22, 24, 25, 32, 35 and Recital 71 GDPR [16]
AI must not turn human beings into objects	<ul style="list-style-type: none"> <li>• Automated decision-making with legal consequences for individuals must be used in a limited scope with appropriate safeguarding measures in place (see TOMs)</li> <li>• Intervention into the automated decision-making process: individuals have the right to request a human intervention (see Individuals rights)</li> <li>• Ability to provide explanation of solely automated decision after it is been made</li> </ul>
AI may only be used for legitimate purposes and may not abrogate the requirement of purpose limitation	<ul style="list-style-type: none"> <li>• PII may only be used for the purpose communicated to and limited to the data acquired from the individuals</li> <li>• Extended purposes must be closely associated with the original purpose</li> <li>• Specifications must include information on the usage of individuals PII for train and/or test data</li> </ul>

<i>Areas and controls</i>	<i>Specification</i>
(see purpose limitation)	<ul style="list-style-type: none"> <li>○ Specify which data of the individual will be used for train/test purposes</li> <li>○ Specify bias mitigation measures incorporated in the model                             <ul style="list-style-type: none"> <li>▪ State the means by which you ensure that the data is representational</li> </ul> </li> </ul>
AI must be transparent, comprehensible and explainable	<ul style="list-style-type: none"> <li>• Transparency of processing is associated with the ease of understanding of the processing activity. It is not enough to explain the result, but rather the end-to-end processes and the decisions made that lead to the result.</li> </ul>
AI must avoid discrimination	<ul style="list-style-type: none"> <li>• Data input sources and data quality must be consistently evaluated to ensure that the principle of fairness, the processing according to the legitimate purpose and the adequacy of the processing is in pace</li> <li>• DPIA results should be evaluated prior to data processing</li> <li>• If data on individuals consist of outlier and model has not been trained such data, incorrect predictions might take place</li> </ul>
The principle of data minimisation applies to AI	<ul style="list-style-type: none"> <li>• Demonstration that the PII is necessary for AI (train/test) purposes and proving the effects to privacy and accuracy if data is not used (e.g., there is no need to collect health information, if the purpose of the fraud detection model is the identification of money laundering transactions)</li> </ul>
AI needs responsibility	<ul style="list-style-type: none"> <li>• Obligation of the controller to demonstrate accountability end-to-end</li> <li>• Ensuring data subject rights</li> <li>• Security and controllability of processing</li> <li>• Conduct DPIA</li> </ul>
AI requires technical and organizational measures (TOMs)	<ul style="list-style-type: none"> <li>• TOMs must be defined for the end-to-end protection of individuals, as the processing of large amounts of data does not dilute the identity of individuals (see Technical and Organizational Measures for further details)</li> </ul>
<b>Individuals rights</b>	Pursuant to Art. 12 - 23 GDPR [16]
to be informed	<ul style="list-style-type: none"> <li>• Informing individuals about the usage of their data for fraud detection purposes supported by AI models before data processing begins</li> <li>• If data is not obtained directly from individuals, they must be notified within one month at latest accordance with Art. 14 GDPR</li> </ul>
of access	<ul style="list-style-type: none"> <li>• Providing individuals access to their data in accordance with Art. 15 GDPR [16]</li> </ul>
to rectification	<ul style="list-style-type: none"> <li>• Wrong data on individuals must be rectified, this is applicable to data stored in the database/raw data and for the pre-processed training data</li> <li>• Attention: a wrong date (e.g., age 32 instead of 23) is not likely to affect the model performance. Nonetheless, the right of an individual must not be disregarded</li> </ul>
to erasure	<ul style="list-style-type: none"> <li>• Erasure of PII in any data processing system, including training datasets</li> </ul> <p>Remember: Erasure of one or few individual’s</p>

Areas and controls	Specification
	data is unlikely to affect the models performance
to restrict processing	<ul style="list-style-type: none"> <li>Individuals have the right to restrict processing of their PII</li> <li>If automated decision making is involved, be able to provide information on how human intervention can replace fully automated decision making procedures</li> </ul>
to data portability	<ul style="list-style-type: none"> <li>Individuals have the right to request their original data in a machine-readable form</li> <li>Data that has been modified in the pre-processing phase may count as PII, but is not affected by the portability request</li> </ul>
to object	<ul style="list-style-type: none"> <li>Individuals can object the usage of their data for AI purposes</li> <li>This may impact their rights to use engage with data processing entities (see Assessment of lawful basis)</li> </ul>
<b>Trade-offs</b>	
Data privacy compliance vs. model accuracy	<ul style="list-style-type: none"> <li>identify and assess any existing or potential trade-offs, when designing or procuring an AI system, and assess the impact it may have on individuals</li> <li>consider available technical approaches to minimize the need for any trade-offs</li> <li>consider any techniques which you can implement with a reasonable level of investment and effort</li> <li>have clear criteria and lines of accountability about the final trade-off decisions. This should include a robust, risk-based and independent approval process</li> <li>Accuracy with respect to privacy: demonstrating the correctness and consistency of personal data</li> <li>Accuracy with respect to statistics: predicting the correct answer; high statistical accuracy (high probability) of predicting the correct answer</li> <li>AI system demonstrate compliance with the fairness principle: higher prediction accuracy means data is PII the                             <ul style="list-style-type: none"> <li>GDPR requires maintenance of correct data (see right to rectification) [8]</li> </ul> </li> <li>Be aware of discrimination:                             <ul style="list-style-type: none"> <li>If model discriminates minorities, due to lack of data about a subject group (e.g., less data on fraud cases associated with woman from particular countries are in dataset; model may lead to wrong prediction)</li> <li>Statistical accuracy (prediction quality) may be increased by feeding more data on minority cases, but may impose higher risks to their privacy, due to additional data</li> <li>A clear process of weighting the interests of privacy rights and statistical accuracy must be defined to mitigate risks</li> </ul> </li> </ul>
<b>Rights in relation to automated decision making and profiling</b>	Pursuant to Art. 22 GDPR [16]

Areas and controls	Specification
Automated decision making models	<ul style="list-style-type: none"> <li>Explanation, if model will entirely make a decision in the respective fraud detection process</li> <li>Demonstrating the ability to provide human intervention on case-by-case basis</li> <li>Demonstrating transparency on the underlying data (see Data minimization and purpose limitation)</li> </ul>
<b>Security and cyber</b>	
Technical and organizational measures (TOMs) [4]	<ul style="list-style-type: none"> <li>Access control: e.g., Access to server rooms only with key or chip card, office rooms secured with alarm</li> <li>Integrity: e.g., user authorizations are restricted to tasks (marketing department only newsletter, accounting also HR data)</li> <li>Pseudonymization: e.g., Replacement of user-related data by random codes</li> <li>Encryption: e.g., Hard disk encryption or cloud solution with encryption</li> <li>Transmission control: e.g., SSL certificate for websites (https://) to transfer data within forms</li> <li>Confidentiality: e.g., password policies</li> <li>Recoverability: e.g., backups that are regularly checked for successful recovery</li> <li>Evaluation: e.g., annual review of technical and organisational measures on effectiveness and plausibility</li> </ul>
Cross-border data transfer security	<ul style="list-style-type: none"> <li>Appropriate safeguards for data stored outside the EU or in transmission from/to EU must ensure that data is pseudonymized, through                             <ul style="list-style-type: none"> <li>encryption when stored</li> <li>encryption in transmission (e.g., transfer takes place via API)</li> <li>TOMs implemented on both ends, i.e., data storage and consumption location</li> <li>Identity and access management guidelines</li> </ul> </li> </ul>
<b>Risk management</b>	
Risk appetite	<ul style="list-style-type: none"> <li>Understanding, evaluating and documenting                             <ul style="list-style-type: none"> <li>risks arising from usage of AI models for the respective processing activity (e.g., risks of ANN for evaluation of new account sign-ups and promo code abuse)</li> <li>risks for the rights and freedom of individual (likelihood and severity)</li> </ul> </li> <li>Mitigation of above mentioned risks</li> </ul>
Special data categories	<ul style="list-style-type: none"> <li>State in a transparent manner how the usage of following data in an AI system will not discriminate individuals                             <ul style="list-style-type: none"> <li>Age</li> <li>Disability</li> <li>Gender reassignment</li> <li>Marriage and civil partnership</li> <li>Pregnancy and maternity</li> <li>Race</li> <li>Religion and belief</li> <li>Sex</li> <li>Sexual orientation</li> </ul> </li> </ul>