

## The Data Minimization Principle

General Data Protection Regulation (GDPR), Europe	gdpr-info.eu/
<p><b>Article 4(1):</b> "personal data" means any information relating to an <b>identified or identifiable natural person</b> ("data subject") [...];</p> <p><b>Article 5(1)(b):</b> Personal data shall be collected for <b>specified, explicit and legitimate purposes and not further processed in a manner that is incompatible with those purposes</b>;</p> <p><b>Article 5(1)(c):</b> Personal data shall be <b>adequate, relevant and limited to what is necessary in relation to the purposes</b> for which they are processed.</p>	
California Privacy Rights Act (CPRA), USA	coppa.ca.gov/
<p><b>Section 1798.100 (a)(1) &amp; (a)(2):</b> [...] A business shall <b>not collect additional categories of (sensitive) personal information</b> or use (sensitive) personal information collected for additional purposes that are incompatible with the <b>disclosed purpose for which the (sensitive) personal information was collected</b> without providing the consumer with notice consistent with this section.</p> <p><b>Section 1798.140 (v)(1):</b> "Personal information" means information that <b>identifies, relates to, describes, is reasonably capable of being associated with, or could reasonably be linked, directly or indirectly, with a particular consumer or household</b>.</p>	
General Personal Data Protection Law (LGPD), Brazil	lgpd-brazil.info/
<p><b>Article 5(l):</b> personal data: information regarding an <b>identified or identifiable natural person</b> ;</p> <p><b>Article 6:</b> Activities of processing of personal data shall be subject to the following principles,</p> <p>I: processing done for <b>legitimate, specific and explicit purposes of which the data subject is informed, with no possibility of subsequent processing that is incompatible with these purposes</b>;</p> <p>II: <b>compatibility of the processing with the purposes communicated to the data subject</b>, in accordance with the context of the processing;</p> <p>III: limitation of the processing to the <b>minimum necessary to achieve its purposes</b>, covering data that are <b>relevant, proportional and non-excessive in relation to purposes of the data processing</b> [...];</p>	
Protection of Personal Information Act (POPIA), South Africa	popia.co.za/
<p><b>Section 1:</b> "personal information" means information relating to an <b>identifiable, living, natural person</b>, [...]</p> <p><b>Section 10:</b> Personal information may only be processed if, given the purpose of processing, it is <b>adequate, relevant and not excessive</b></p> <p><b>Section 13(1):</b> Personal information must be collected for a <b>specific, explicitly defined and lawful purpose</b> [...]</p>	
Consumer Data Rights (CDR), Australia	www.legislation.gov.au/Details/F2023C00735
<p><b>Rule 1.8(a):</b> An accredited person complies with the data minimisation principle if when making a consumer data request on behalf of a CDR consumer, it does not seek to collect: (i) <b>more CDR data than is reasonably needed</b>; [...]</p> <p><b>Rule 4.11(1)(a):</b> When asking a CDR consumer to give a consent, an accredited person must allow the CDR consumer to choose the types of CDR data to which the consent will apply by enabling the CDR consumer to actively select or otherwise clearly indicate: (i).; and (ii) in the case of a use consent – <b>the specific uses of collected data to which they are consenting</b>; [...]</p>	
Personal Information Protection Act (PIPA), South Korea	www.pipc.go.kr/eng/index.do
<p><b>Article 2(1):</b> The term "personal information" means any of the following information relating to a living individual: (a) Information that <b>identifies a particular individual</b> [...]; (b) Information which, even if it by itself does not identify a particular individual, <b>may be easily combined with other information to identify a particular individual</b> [...];</p> <p><b>Article 3(1):</b> The personal information controller shall <b>specify explicitly the purposes for which personal information is processed</b>; and shall collect personal information <b>lawfully and fairly to the minimum extent necessary for such purposes</b>.</p>	

## How do we formalize data minimization in ML?

## Data Minimization as Optimization

The goal of *data minimization* can be modelled as minimizing the size of the dataset  $D$  by removing entries from the features  $\mathbf{x}_i$  (**data relevance**) while also retaining similar performance as it would have on the complete dataset (**purpose**).

$$\underset{\mathbf{B} \in \{0,1\}^{n \times p}}{\text{Minimize}} \quad \|\mathbf{B}\|_1 \quad \text{s.t. :} \quad J(\hat{\theta}; \mathbf{X}, \mathbf{Y}) - J(\emptyset; \mathbf{X}, \mathbf{Y}) \leq \alpha$$

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n \ell(f_{\theta}(\mathbf{x}_i \odot \mathbf{B}_i), y_i)$$

### Purpose limitation

Data collectors must define a clear, legal objective before data collection and use the data solely for that objective. In ML, this *purpose* translates to collecting data solely for training models to achieve optimal performance, for a clearly defined performance metric.

### Data relevance

Data minimization aims to remove data that does not serve the purpose. In ML contexts, this translates to retaining only data that contributes to the performance of the model. We interpret this as an optimization, minimizing data while maintaining performance.

### Data Privacy

Data minimization aims to avoid the collection of any unnecessary information that could compromise an individual's privacy, thus creating an implicit expectation of privacy from minimization.

#### EU AI Act: Recital 69

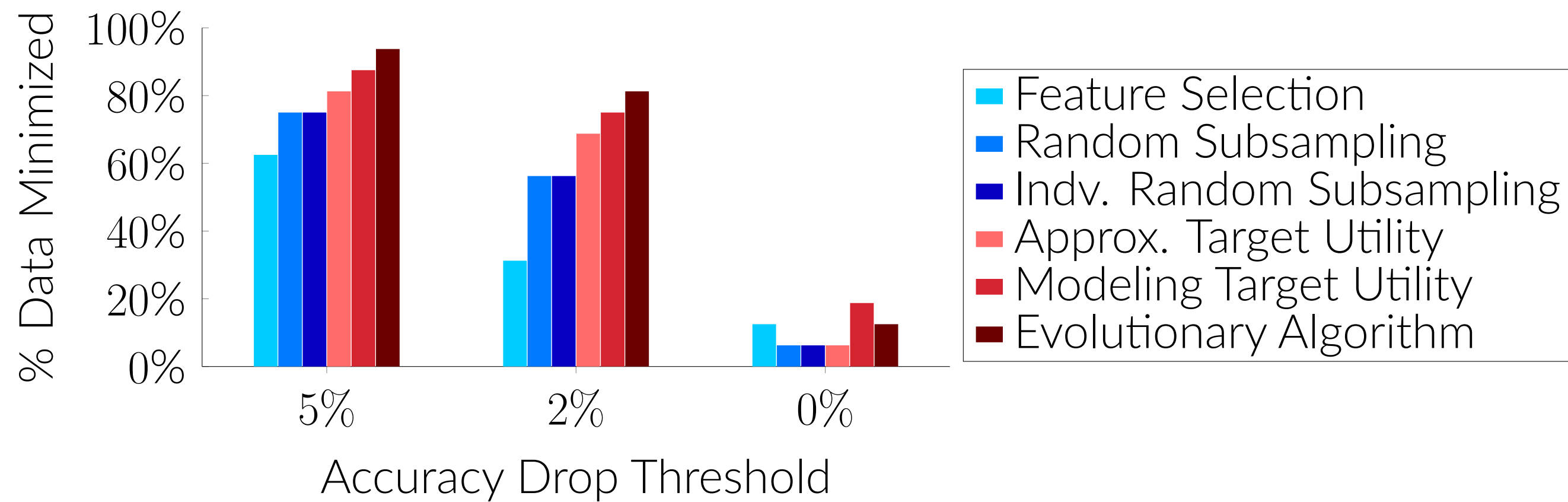
The right to privacy and to protection of personal data must be guaranteed throughout the entire lifecycle of the AI system. In this regard, the principles of data minimisation and data protection by design and by default, as set out in Union data protection law, are applicable when personal data are processed.

*Do current data minimization requirements in regulations genuinely meet privacy expectations?* **Unfortunately, No!**

## Empirical Analysis

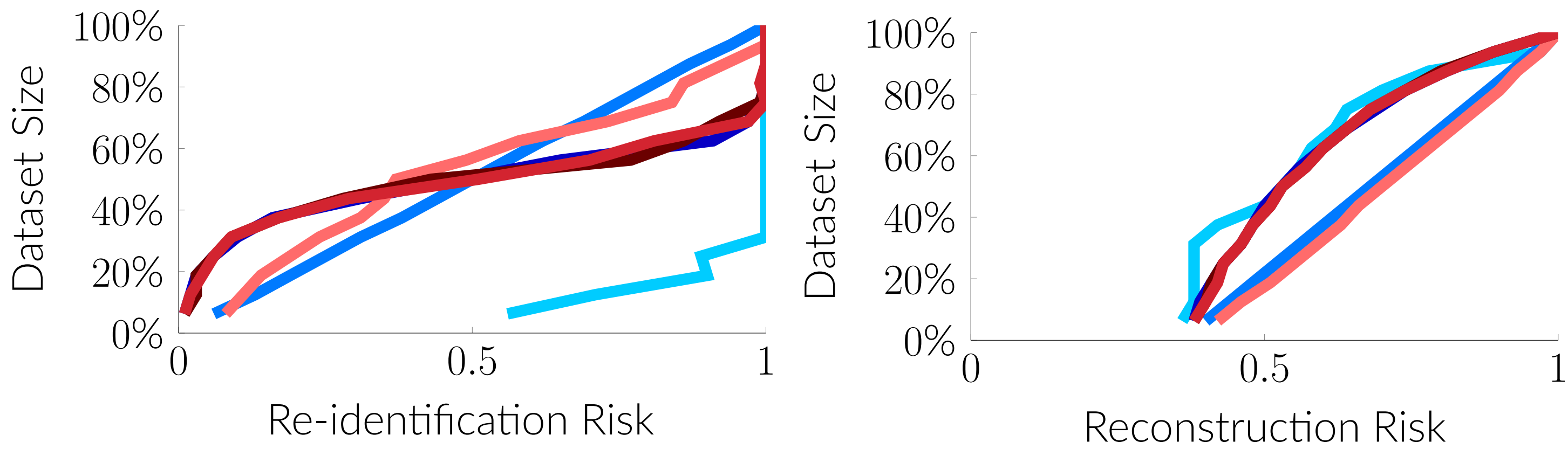
### Utility

A substantial amount of data can often be removed without sacrificing utility, suggesting that much of the collected data is superfluous in the datasets analyzed.



### Privacy

There exists a fundamental misalignment between the goals of data minimization and actual privacy.



## Open Challenges

- How to more efficiently solve the optimization challenges that come with minimization?
- Can other formalisms be created to better operationalize data minimization that incorporates nuances like active minimization during data collection, instead of data removal?
- What ethical and fairness considerations of data minimization are needed to make sure we don't leave the minority more vulnerable to privacy risks?