

## Table of De-Identification Techniques

Name of Technique	Description / Examples	Pros	Cons
<b>Redaction</b>	<b>Erasing or expunging</b> sensitive data from a record.	Reduces risk if data are disclosed inadvertently or through unauthorized access; useful when the erased data elements are not needed for analysis (typical with direct identifiers).	Not effective if done improperly (e.g., if the erasure can be reversed or if enough indirect identifiers remain).
<b>Suppression</b>	<p><b>Removing</b> data (e.g., from a cell or row in a table, or data element(s) in a record) prior to dissemination to prevent the identification of individuals in small groups or those with unique characteristics.</p> <p>Examples:</p> <ul style="list-style-type: none"> <li>• Suppressing the value of a single field, such as a field in a patient record containing a very rare disease.</li> <li>• Not reporting observations for those patients where the number of patients for any combination of zip code, age, and diagnosis is below a given threshold (e.g., 5 people).</li> </ul>	<p>Useful when multiple indirect identifiers pose a risk for re-identification.</p> <p>More easily done with tabular data.</p> <p>Helpful when presenting analysis of findings to the institution that provided the data.</p> <p>Helpful in public health reporting.</p>	<p>May result in minimal data being produced for small populations, and it usually requires additional suppression of non-sensitive data to ensure adequate protection of PII (e.g., complementary suppression of one or more non-sensitive cells in a table so that the values of the suppressed cells may not be calculated by subtracting the reported values from the row and column totals).</p> <p>Can be difficult to perform properly.</p> <p>Is less likely to be effective if there are additional data available elsewhere.</p>
<p><b>Blurring:</b></p> <ul style="list-style-type: none"> <li>• Aggregation</li> <li>• Generalization</li> <li>• Pixelation</li> </ul>	<p><b>Reducing precision</b> of data by combining one or more data elements.</p> <p><b>Aggregation:</b> combining individual subject data with a sufficient number of other subjects to disguise the attributes of a single subject (e.g., reporting a group</p>	<p>Minimizes risk of identification by focusing on collective data rather than individual data.</p> <p>Useful for “big picture” analyses.</p>	<p>Decreases reliability of data and increases potential for false conclusions.</p> <p>Aggregation: may not be possible with a small pool of subjects.</p>

Name of Technique	Description / Examples	Pros	Cons
	<p>average instead of an individual value).</p> <p><b>Generalization:</b> collecting or reporting values in a given range (e.g., using age or age-range instead of date of birth); including individual data as a member of a set (e.g., creating categories that incorporate unique cases); or reporting rounded values instead of exact amounts.</p> <p><b>Pixelation:</b> modifying or obscuring visual information (e.g., blurring out faces in a photograph).</p>		<p>Generalization: unhelpful for case studies or in situations where details or specificity enhance findings.</p> <p>Pixelation: technology exists to reverse such modifications; other factors in a photo can lead to re-identification, such as setting and clothing.</p>
<p><b>Masking:</b></p> <ul style="list-style-type: none"> <li>• Pseudonymization</li> <li>• Coding</li> <li>• Perturbation</li> <li>• Randomization</li> <li>• Swapping</li> <li>• Shuffling</li> <li>• Scrambling</li> <li>• Encryption</li> <li>• Noise</li> <li>• Differential Privacy</li> </ul>	<p><b>Replacing</b> one data element with either a random or made-up value, or with another value in the data set; can be done manually or by using an algorithm.</p> <p><b>Pseudonymization/Coding:</b> replacing a real name with a made up name or a real value with a made-up value.</p> <p><b>Perturbation:</b> replacing sensitive info with realistic but inauthentic data or modifying original data based on predetermined masking rules (which may include randomization). Example: an algorithm which replaces the date of birth of subjects.</p> <p><b>Swapping/Shuffling:</b> data for one or more variables are switched with another record, so that the data user</p>	<p>Attempts to retain the functional usability of the data while concealing information that could lead to identification.</p> <p>Pseudonymization/Coding: allows for a unique descriptor to trace data across multiple records; useful for multiple data instruments.</p> <p>Perturbation: reduces the likelihood of reverse identification.</p> <p>Swapping/Shuffling: useful for creating data sets for software testing where fields must be</p>	<p>Can decrease accuracy of computations in some cases, affecting validity of data.</p> <p>Techniques may be ineffective for small data sets.</p> <p>Algorithms used for masking can be reverse engineered.</p>

Name of Technique	Description / Examples	Pros	Cons
	<p>does not know whether the real data values correspond to certain records (i.e., all the values in the data set are real, but are assigned to the wrong people).</p> <p><b>Scrambling/Encryption:</b> data are algorithmically scrambled and only those with access to the appropriate key can view the encrypted data.</p> <p><b>Noise/Differential Privacy:</b> statistical technique that introduces errors by randomly misclassifying values of categorical variable(s).</p>	<p>present and have realistic looking values.</p> <p>Noise/Differential Privacy: allows for quantification of potential privacy loss, enabling a more accurate risk assessment; useful for large data sets.</p>	
<b>Subsampling</b>	<b>Releasing either a representative or random subsample</b> of data instead of an entire data set.	Minimizes risk of identification by reducing the amount of data reported.	May not yield representative and generalizable estimates of a study's overall subject population.