

Presentation Data Redaction

Intro

"Data anonymization" is the process of removing or altering personal information that can identify a specific individual from a dataset to protect their privacy. This makes it impossible to trace the information back to any one person. Personal information includes data such as ID numbers, tax IDs, patient numbers, bank account or credit card numbers, home addresses, phone numbers, and even biometric data like photos, X-rays, fingerprints, eye scans, voice signatures, and facial recognition data. It also includes other details that, when combined with information such as birthdate, birthplace, race, religion, location data, job history, medical records, or education, could still be used to identify an individual.



Part I

These days, everyone is at risk of their data getting leaked or being hit by a cyberattack, no matter how big or small the company is. When a data breach happens, sensitive information can get stolen and sold on the dark web or to other parties, just like what happened with the Microsoft data leak back in 2021. When it comes to keeping machine transactions secure, a key part of protecting data is anonymization.

For example, an e-commerce company like Amazon could use data anonymization to keep their customers' privacy safe when they're looking at shopping habits. When they analyze transaction data, they might swap out customer names for unique IDs, like 'Customer12345'. This helps them figure out shopping trends and preferences – like which products are often bought together – without revealing who the actual customers are.

Data anonymization tools are basically software designed to protect this anonymized information by changing those personal identifiers in a dataset. They tweak or hide the data so it's not easy to link the original information back to a specific person, which keeps people's data private. Let's dive a bit deeper into how these tools work.

Part II

Data anonymization tools use a few different tricks to protect data in a way that lines up with GDPR:

- Data Masking: This is where they replace the real data with fake but realistic-looking stuff that doesn't actually mean anything. For instance, a customer's name like 'John Smith' could become 'Jane Doe'. This way, the data is still useful, but the person's identity is safe. You can read more about what data masking is on the Lingvanex blog.
- -Pseudonymization: This method swaps out identifying info with nicknames or new values. So, a user's email address might get replaced with a unique code like 'user123'. You can usually only reverse this with a special key, which means the original data can be safely recovered if needed.
- Generalization: This technique makes the data less specific to protect privacy. For example, instead of saying someone is exactly 29 years old, the data might show an age range like '25-30'. This helps keep things private while still giving useful info for analysis.
- Perturbation: This involves adding a little 'noise' to the data, slightly messing it up so you can't pinpoint exact individuals, but the overall trends still hold true. For example, if sales data shows 100 units of something were sold, perturbation might change that to 98 or 102. This keeps the data useful for analysis without revealing precise figures.

By doing these things, organizations can analyze anonymized information while still protecting

Practical work

Red and black colors are often used in website design to create contrast and attract attention. Red symbolizes energy, passion and urgency, while black symbolizes elegance, strength and mystery. Their combination can be very effective for accentuating important elements such as call-to-action buttons or important messages. That's why I decided to use these colors for website design, as well as the logo.

Well I created the header rather simply, it doesnt contain much, but the 2 relativly hard parts are the export and account button. One creates a panel that contains 3 buttons, export all, clear all and close. the mechanism behind export all is that, whenever a user uses the service, for example scans a pdf, that scanned redacted pdf gets saved in the browser even if the user didnt save it themselves it uses IndexedDB, at first i though of using the local storage of a browser but ran into problems when the size of everything reached 5MB, so I had to do the harder way. The account button summons another panel, with sign in with key, sign up. after signing in or signing up we are greeted with 3 buttons, download login key, export all images, and sign out. you can use the login key to sign in later and keep ur pass and username, the sign out button prompts you to save before exiting and wipes the data stored in the indexedDB.

Data anonymization tools use a few different tricks to protect data in a way that lines up with GDPR: Moving on to the Hero, there is nothing to say, no scripts were used, only html and some css.

We then arrive at the Text Anonymization section where the actual magic starts, the section contains 2 main parts, 1 is a legacy version of the "task" and a beta version. The legacy uses simple algorithms and the beta uses compromise.js, it correctly censors most western names, and has a bonus function to "scramble" the answers, so that the text wont just become "REDACTED" and all that but also change names, cards, numbers and etc into fake ones using https://cdn.sky-pack.dev/@faker-js/faker. and well there is a download button and a redacted info count, not much we can say about that.

The next section is image redaction section, you upload nearly any image file, it scans it, and censors the face, but it doesnt censor words, for that you will have to put the img into a pdf. There exists a download button and also a blur intensity slider.

The hard part started here, PDF scanning, finding text, finding location of redacted words, and making then blacked out, and then blurring the faces. The word location took me hours to do because I decided to use the already "redacted" image as a location to put my blacked out squares instead of the raw PDF. there is also a bonus feature, the script find the faces, and cuts them out so you can download them seperatly.

We find ourselves in the second last section, the chat system, it basically combines all the previosuly mentioned redaction methods into a fake chat, it doesnt use AI so it doesnt "leak" anything into some other server.

The footer doesnt have much to say about it, it doesnt contain scripts and is a basic footer.

Informative overview

here is another small overview:

Technologies Used

Frontend Framework: Pure HTML/CSS/JavaScript (no frameworks)

**Key Libraries:

PDF.js for PDF processing

face-api.js for face detection

Compromise.js for NLP text processing

Faker.js for generating fake data

JSZip for creating zip archives

Storage: IndexedDB for local storage of processed images

Modules: ES6 modules for code organization

=======

Privacy & Security Implementation followed several important rules for a privacy-focused application: **Local Processing:** All processing happens client-side (no server calls) Sensitive data never leaves the user's browser Confirmed by examining all JavaScript files Data Handling: Implemented proper redaction for: Personal names (using NLP) IIN numbers (12-digit numbers) Credit card numbers (16-20 digits) **Email addresses** Provided both redaction (blacking out) and scrambling (replacement with fake data) options **Authentication:** Local account system stores credentials only in localStorage Allows export of account data as JSON Clear warning before sign-out about data loss **GDPR Compliance Indicators:** Hero section badges mention "GDPR Ready" All processing stays local as promised No evidence of analytics or tracking File Handling: Images and PDFs processed entirely in browser Face detection and blurring happens client-side Option to export processed files

Implementation Details

Text Anonymization:

Two modes (legacy pattern matching and smart NLP-based)

Proper handling of Cyrillic and Latin characters

Counts of redacted items displayed

Image Processing:

Face detection with adjustable blur intensity

Canvas-based processing

Download options for processed images

PDF Processing:

Text extraction and redaction

Face detection in PDF pages

Multi-page handling

Download options per page

Chat Interface:

Unified interface for all processing types

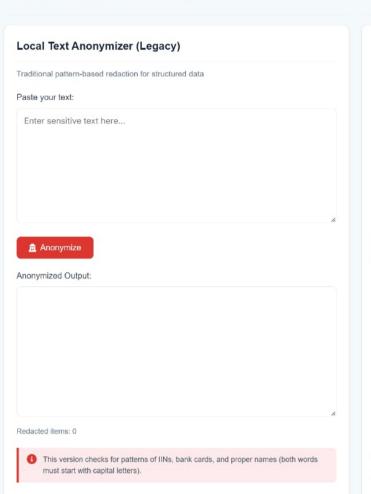
Conversation history maintained

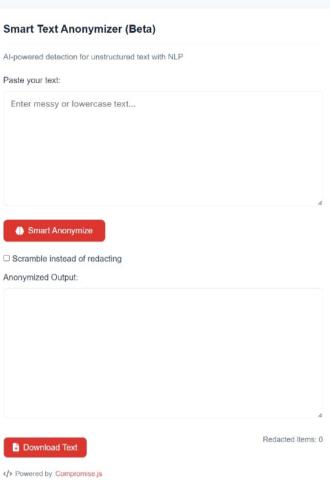
Clear mode switching



Text Anonymization

Protect sensitive information in your text documents

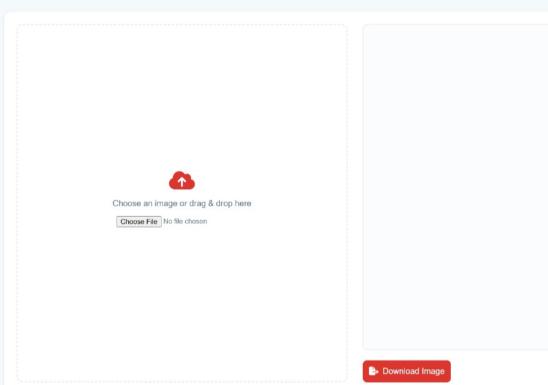




Blur Intensity:

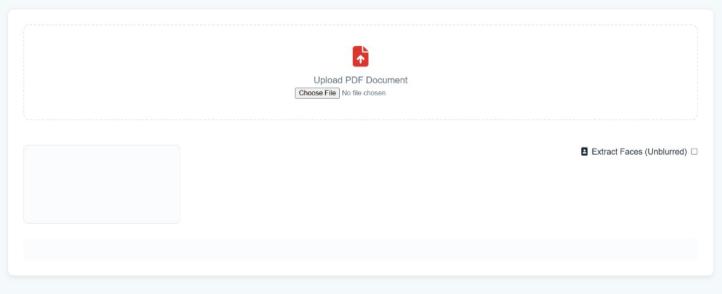
Image Anonymization

Detect and blur faces in images while preserving quality



PDF Anonymization

Redact text and blur faces in PDF documents



Try out the Chat version of the system





Thanks for listening.